# UNVEILING THE COMPLEXITY OF HUMAN MOBILITY BY MINING & QUERYING MASSIVE TRAJECTORY DATA

**Fosca Giannotti**

Knowledge Discovery & Data Mining  LAB ISTI-CNR & Università di Pisa
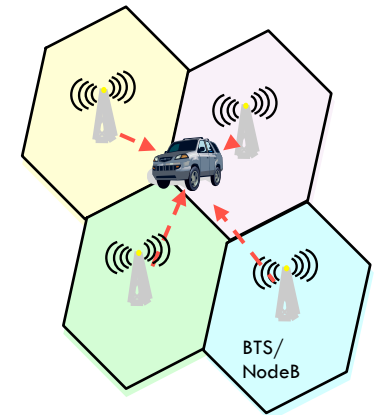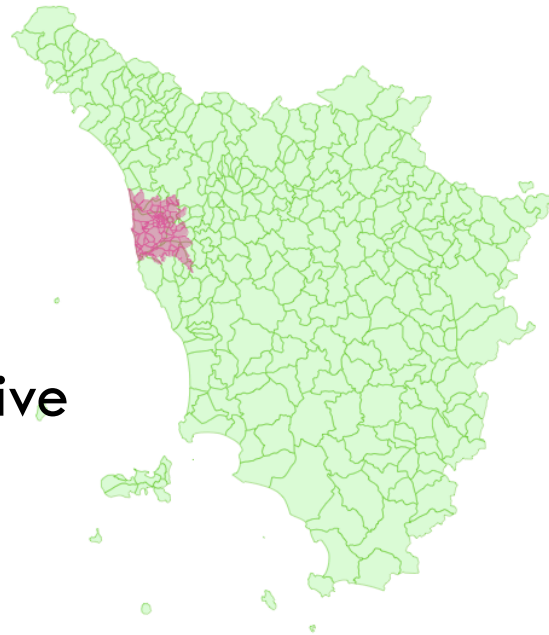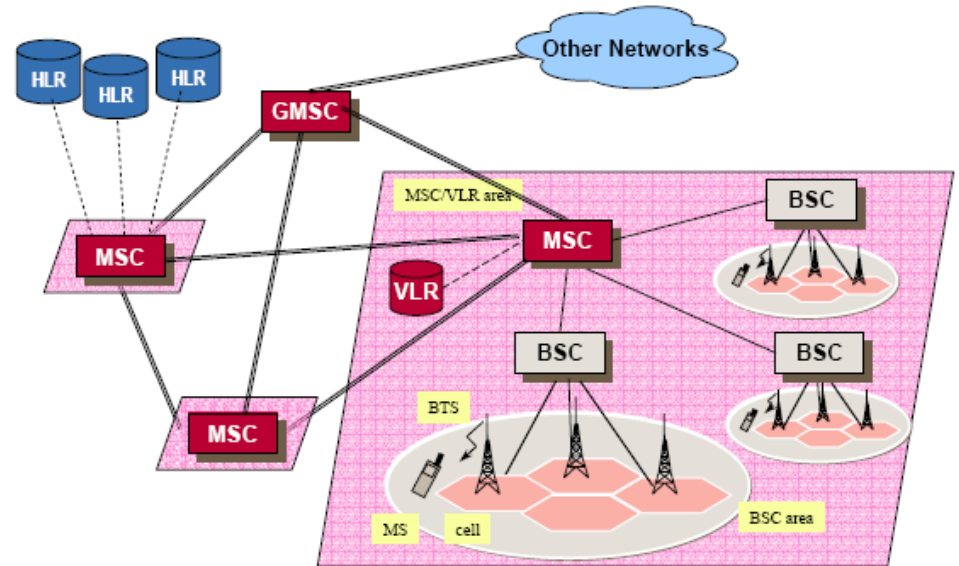
**http://kdd.isti.cnr.it**

# BIG DATA as a proxy of human mobility

# GSM data



- Mobile Cellular Networks handle information about the positioning of mobile terminals
  - CDR Call Data Records: call logs (tower position, time, duration,..)
  - Handover data: time of tower transition
- More sophisticated Network Measurement allow tracking of all active (calling) handsets

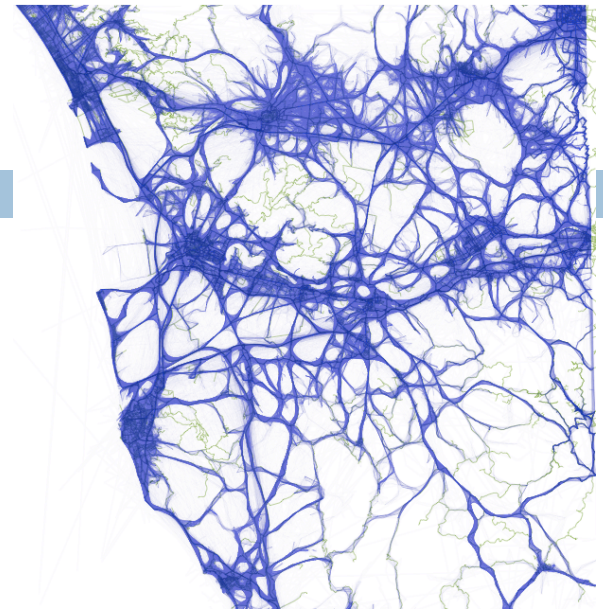# GSM data as a proxy of presence and fluxes

Video: Paris_splines.avi

# GPS tracks



□ Onboard navigation devices send GPS tracks to central servers

  □ Sampling rate ~3 secs
  □ Spatial precision ~ 10 m



Ide;Time;Lat;Lon;Height;Course;Speed;PDOP;State;NSat
…
8;22/03/07 08:51:52;50.777132;7.205580; 67.6;345.4;21.817;3.8;1808;4
8;22/03/07 08:51:56;50.777352;7.205435; 68.4;35.6;14.223;3.8;1808;4
8;22/03/07 08:51:59;50.777415;7.205543; 68.3;112.7;25.298;3.8;1808;4
8;22/03/07 08:52:03;50.777317;7.205877; 68.8;119.8;32.447;3.8;1808;4
8;22/03/07 08:52:06;50.777185;7.206202; 68.1;124.1;30.058;3.8;1808;4
8;22/03/07 08:52:09;50.777057;7.206522; 67.9;117.7;34.003;3.8;1808;4
8;22/03/07 08:52:12;50.776925;7.206858; 66.9;117.5;37.151;3.8;1808;4
8;22/03/07 08:52:15;50.776813;7.207263; 67.0;99.2;39.188;3.8;1808;4
8;22/03/07 08:52:18;50.776780;7.207745; 68.8;90.6;41.170;3.8;1808;4
8;22/03/07 08:52:21;50.776803;7.208262; 71.1;82.0;35.058;3.8;1808;4
8;22/03/07 08:52:24;50.776832;7.208682; 68.6;117.1;11.371;3.8;1808;4
…

# GPS: detailed movements within an area

Video: moves_viz_prov_cut.mov

# GPS: movements within the town

Video: moves_viz_city_cut.mov

# Social Networks: goal of movements

Video: flickr_cut.mov

# Plan of the presentation

- Mastering the overall KDD process
  - M-atlas platform
- Exemplar case studies
  - Advanced OD Matrix browsing
  - Understanding collective patterns
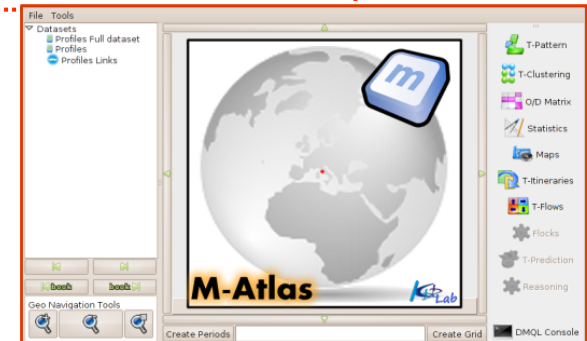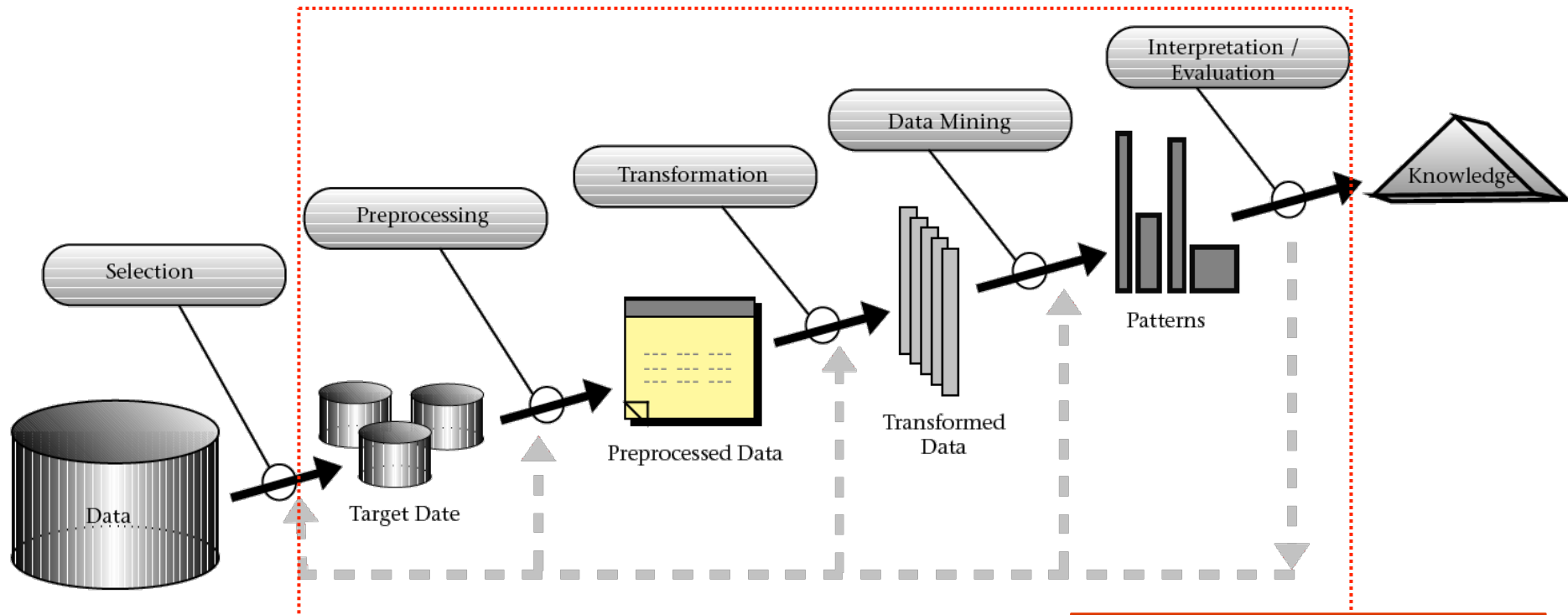  - Understanding Individual profiles
  - Putting interactions in the game

# Mastering the overall KDD process: M-Atlas platform

Fosca Giannotti · Mirco Nanni · Dino Pedreschi · Fabio Pinelli · Chiara Renso · Salvatore Rinzivillo · Roberto Trasarti

Unveiling the complexity of human mobility by querying and mining massive trajectory data
*The VLDB Journal*, 2011

Roberto Trasarti, Fosca Giannotti, Mirco Nanni, Dino Pedreschi, Chiara Renso.
A Query Language for Mobility Data Mining.
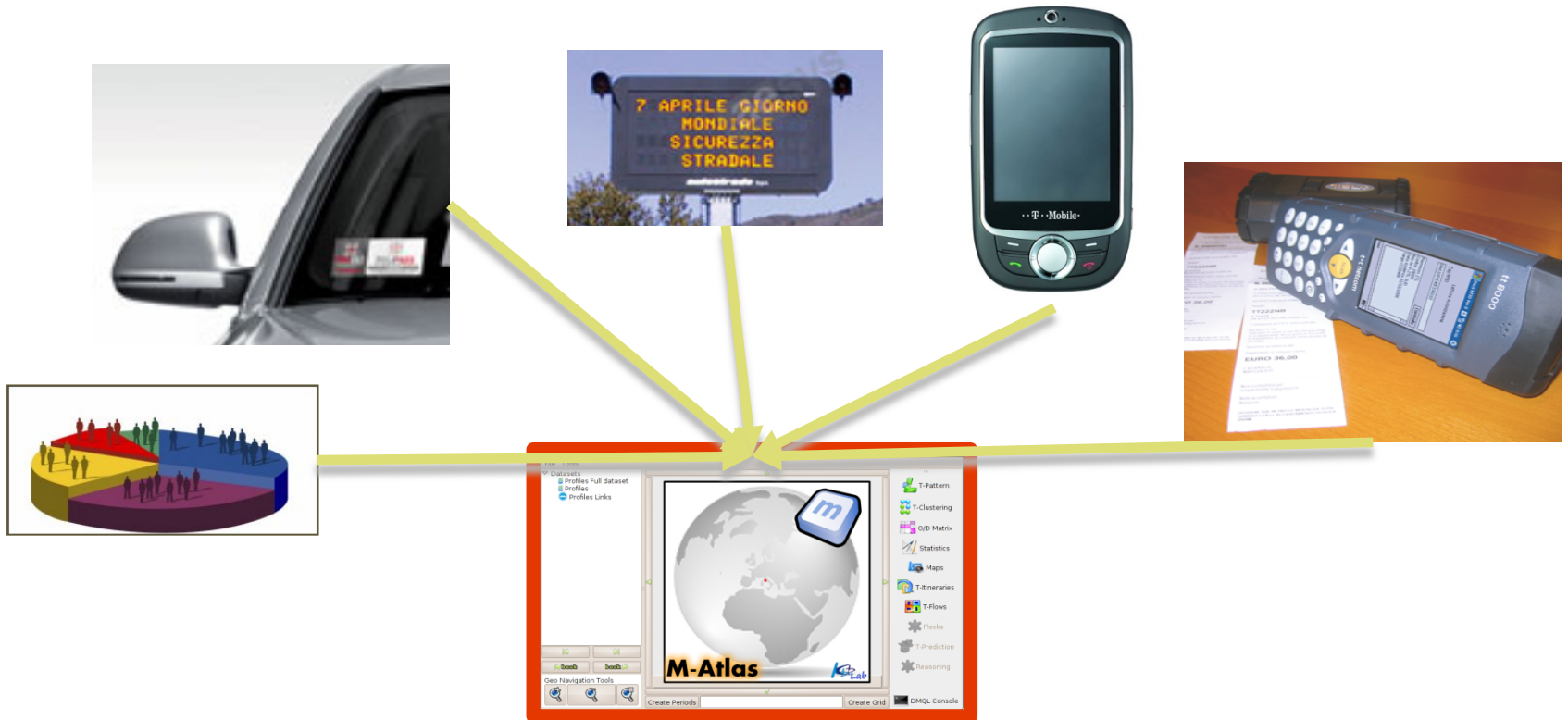*International Journal of Data Warehousing and Mining* (IJDWM) 2010

# Knowledge Discovery process
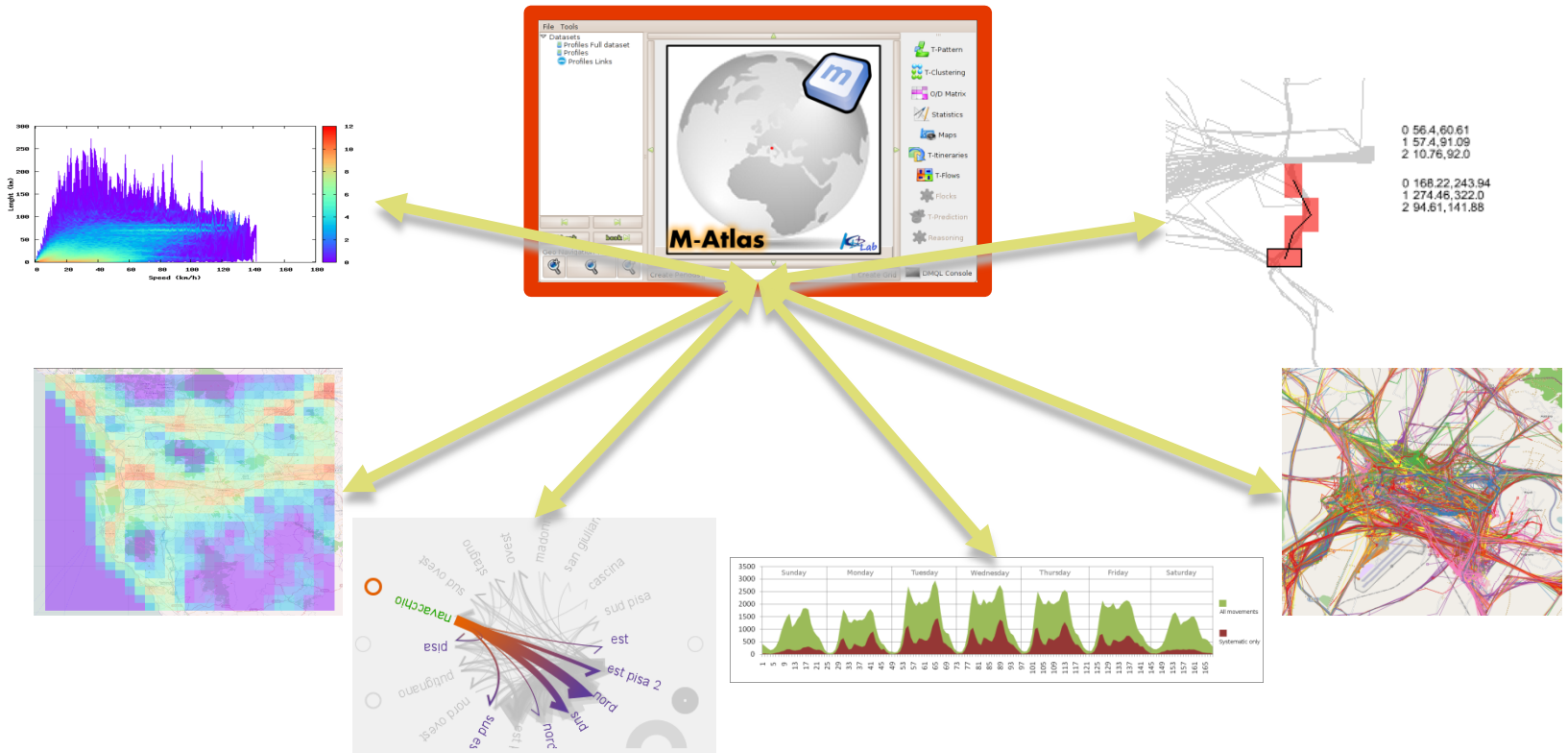
# M-Atlas platform

M-Atlas: An analytical system to create and navigate an atlas of urban mobility

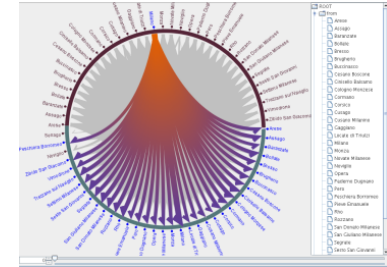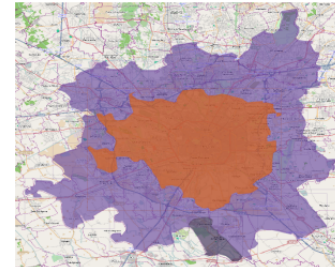Source data: GPS, GSM, Sensors, Rfid, spatial data

# M-Atlas platform

A tool kit to extract, store, combine different kinds of models to build mobility knowledge discovery processes.
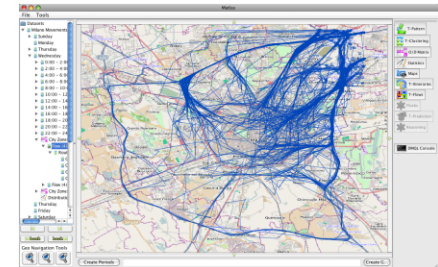
# DMQL EXPRESSIVENESS:
## How do people leave the city toward suburban areas?
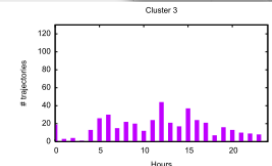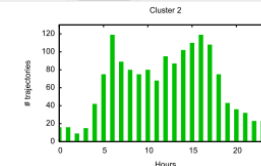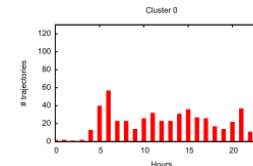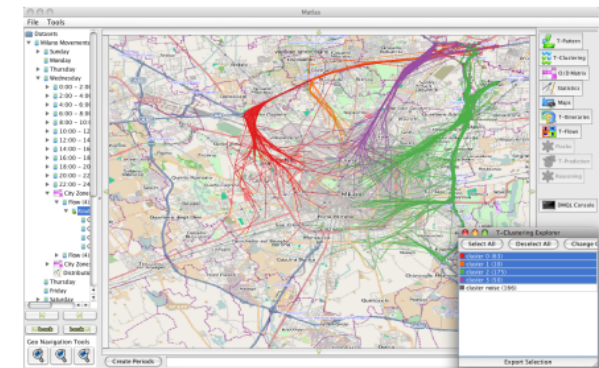
CREATE MODEL MilanODMatrix AS MINE ODMATRIX

FROM (SELECT t.id, t.trajectory FROM TrajectoryTable t),

(SELECT orig.id, orig.area FROM MunicipalityTable orig),

(SELECT dest.id, dest.area FROM MunicipalityTable dest)


CREATE RELATION CenterToNESuburbTrajectories USING ENTAIL

FROM (SELECT t.id, t.trajectory FROM TrajectoryTable t, MilanODMatrix m

WHERE m.origin = Milan AND

m.destination IN (Monza, ..., Brugherio))


CREATE MODEL ClusteringTable AS MINE T-CLUSTERING

FROM (Select t.id, t.trajectory from CenterToNESuburbTrajectories t)

SET T-CLUSTERING.FUNCTION = ROUTE_SIMILARITY AND

T-CLUSTERING.EPS = 400 AND

T-CLUSTERING.MIN_PTS = 5

# A DataWarehouse for OD Matrix

# OD Matrix

- Model mobility demand by measuring the flows among different areas

- General approach
  - Spatial grid with relevant zones
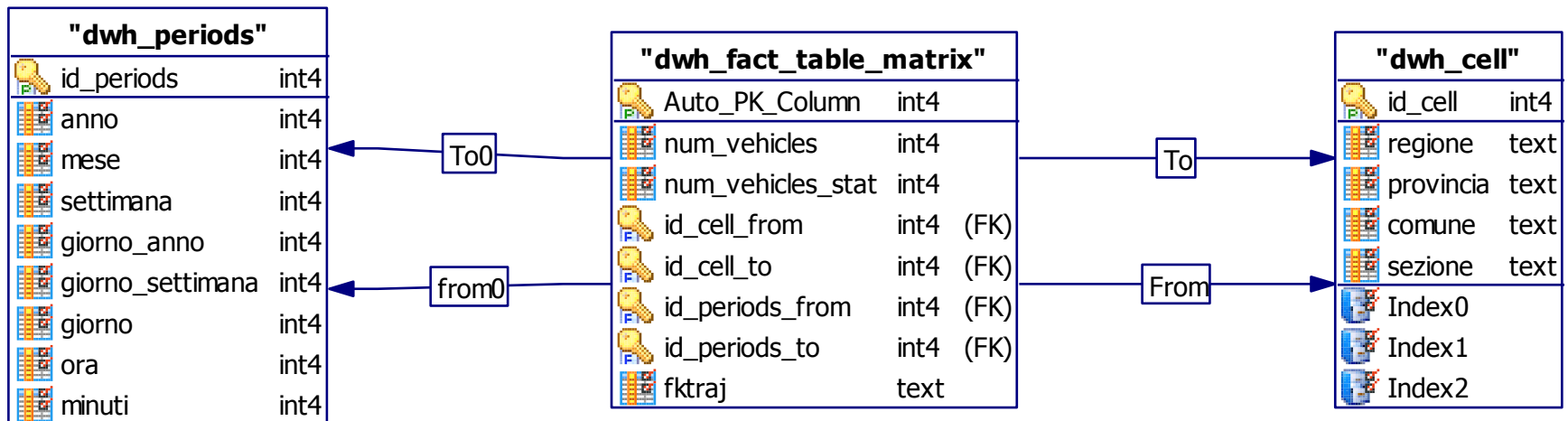  - (Estimated) flows of movement from origin to destination

U Migration

# OD Matrix exploration

- □ OD Matrix should answer the questions
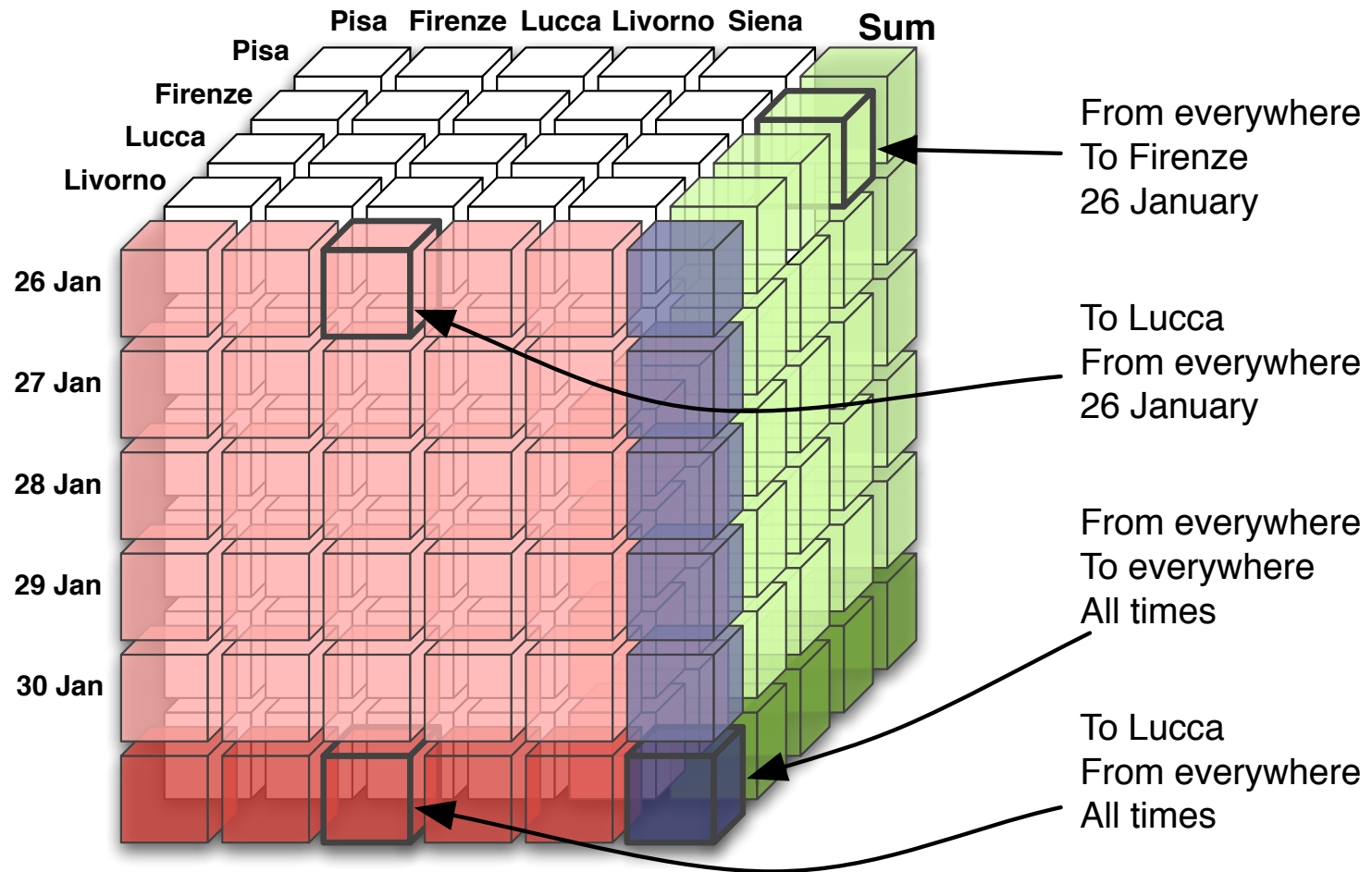  - ▪ From which region?
  - ▪ To which region?
  - ▪ When?
  - ▪ How many?
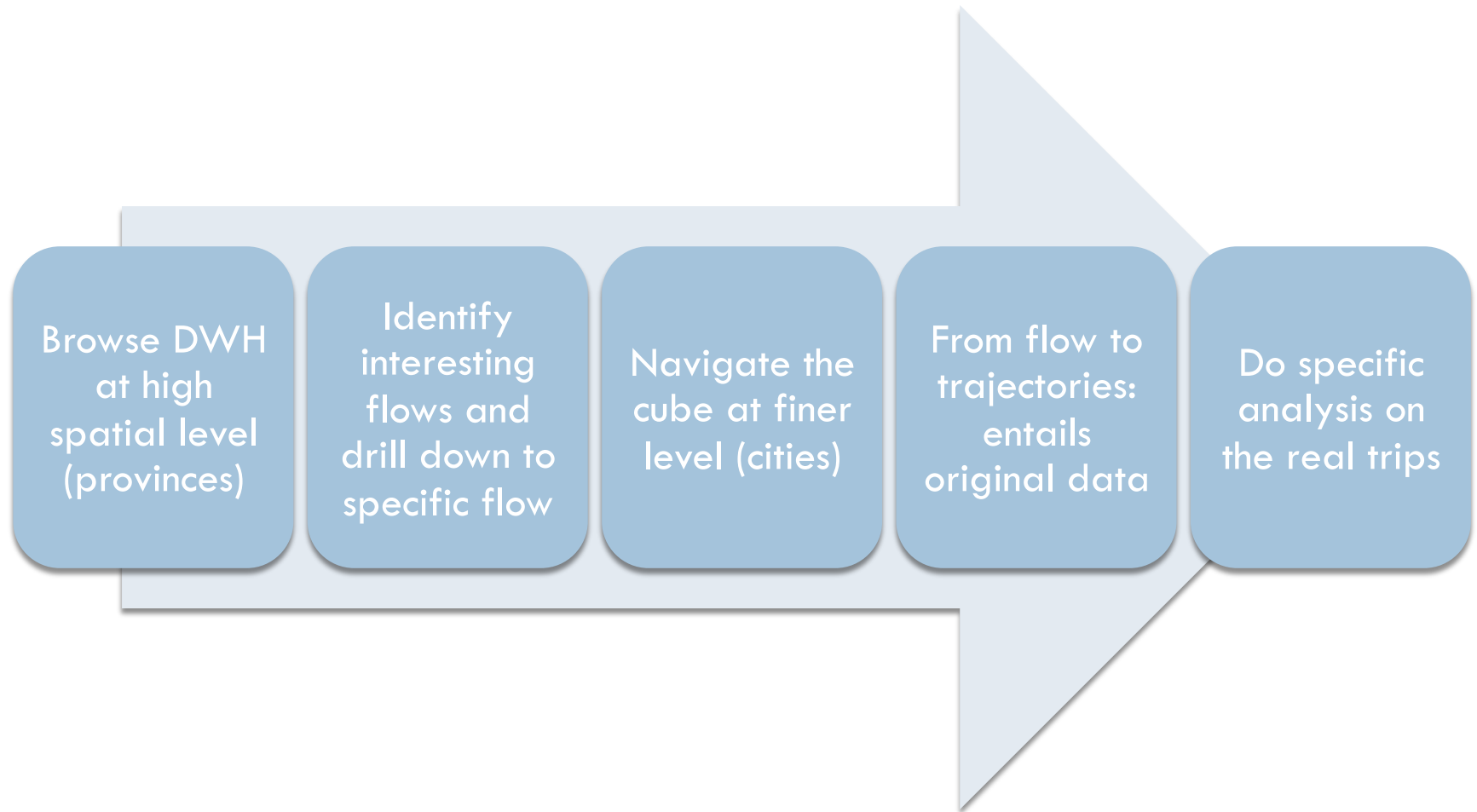
- □ DW Concepts
  - ▪ Facts: basic observation
    - ▪ Aggregated movements from an origin to a destination
  - ▪ Dimensions
    - ▪ Origins
    - ▪ Destinations
    - ▪ Time
  - ▪ Measures
    - ▪ Count
    - ▪ Ratio over total

**"dwh_periods"**

| | |
|---|---|
| id_periods | int4 |
| anno | int4 |
| mese | int4 |
| settimana | int4 |
| giorno_anno | int4 |
| giorno_settimana | int4 |
| giorno | int4 |
| ora | int4 |
| minuti | int4 |

**"dwh_fact_table_matrix"**

| | | |
|---|---|---|
| Auto_PK_Column | int4 | |
| num_vehicles | int4 | |
| num_vehicles_stat | int4 | |
| id_cell_from | int4 | (FK) |
| id_cell_to | int4 | (FK) |
| id_periods_from | int4 | (FK) |
| id_periods_to | int4 | (FK) |
| fktraj | text | |

To0    from0    To    From

**"dwh_cell"**

| | |
|---|---|
| id_cell | int4 |
| regione | text |
| provincia | text |
| comune | text |
| sezione | text |
| Index0 | |
| Index1 | |
| Index2 | |

# OD Matrix: DW design

# The general process

Browse DWH at high spatial level (provinces)

Identify interesting flows and drill down to specific flow

Navigate the cube at finer level (cities)

From flow to trajectories: entails original data

Do specific analysis on the real trips

# Navigate the cube at higher spatial level (provinces): pivot table

| Time to | Cell to | | Cell from | | Measures | |
|---|---|---|---|---|---|---|
| (All) | Regione | Provincia | Regione | Provincia | ▼ Numero Veicoli | ▣ Perc |
| +All Time To | Toscana | +Pisa | Toscana | +Pisa | 462.583 | 86,53% |
| | | | | +Firenze | 27.742 | 13,26% |
| | | | | +Livorno | 20.429 | 10,05% |
| | | | | +Lucca | 17.681 | 04,07% |
| | | | | +Pistoia | 5.727 | 01,30% |
| | | +Pistoia | Toscana | +Pistoia | 405.003 | 92,05% |
| | | | | +Lucca | 19.040 | 04,38% |
| | | | | +Firenze | 7.853 | 03,75% |
| | | | | +Pisa | 5.630 | 01,05% |
| | | | | +Livorno | 2.306 | 01,13% |
| | | +Lucca | Toscana | +Lucca | 388.854 | 89,42% |
| | | | | +Pistoia | 19.268 | 04,38% |
| | | | | +Pisa | 17.750 | 03,32% |
| | | | | +Livorno | 6.488 | 03,19% |
| | | | | +Firenze | 2.747 | 01,31% |
| | | +Firenze | Toscana | +Firenze | 163.845 | 78,34% |
| | | | | +Pisa | 27.571 | 05,16% |
| | | | | +Pistoia | 7.769 | 01,77% |
| | | | | +Livorno | 6.617 | 03,26% |
| | | | | +Lucca | 2.650 | 00,61% |
| | | +Livorno | Toscana | +Livorno | 167.347 | 82,36% |
| | | | | +Pisa | 21.088 | 03,94% |
| | | | | +Firenze | 6.971 | 03,33% |
| | | | | +Lucca | 6.625 | 01,52% |
| | | | | +Pistoia | 2.228 | 00,51% |

- The cube dimensions are flattened by means of a multi-row table

- Example at the province level:
  - How many trips from Lucca province to Pisa province?
  - How many in the other way?

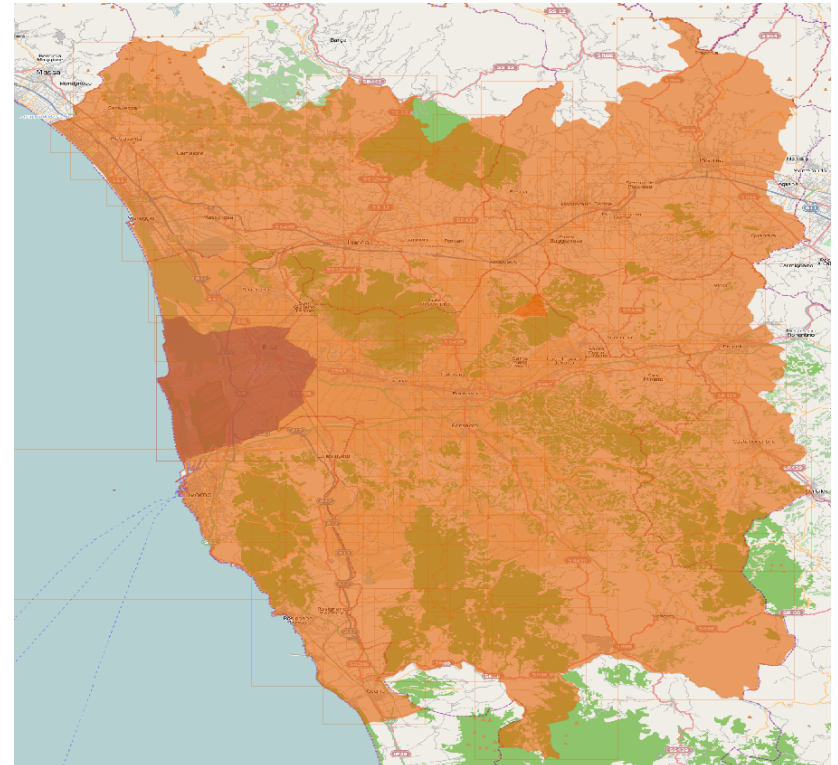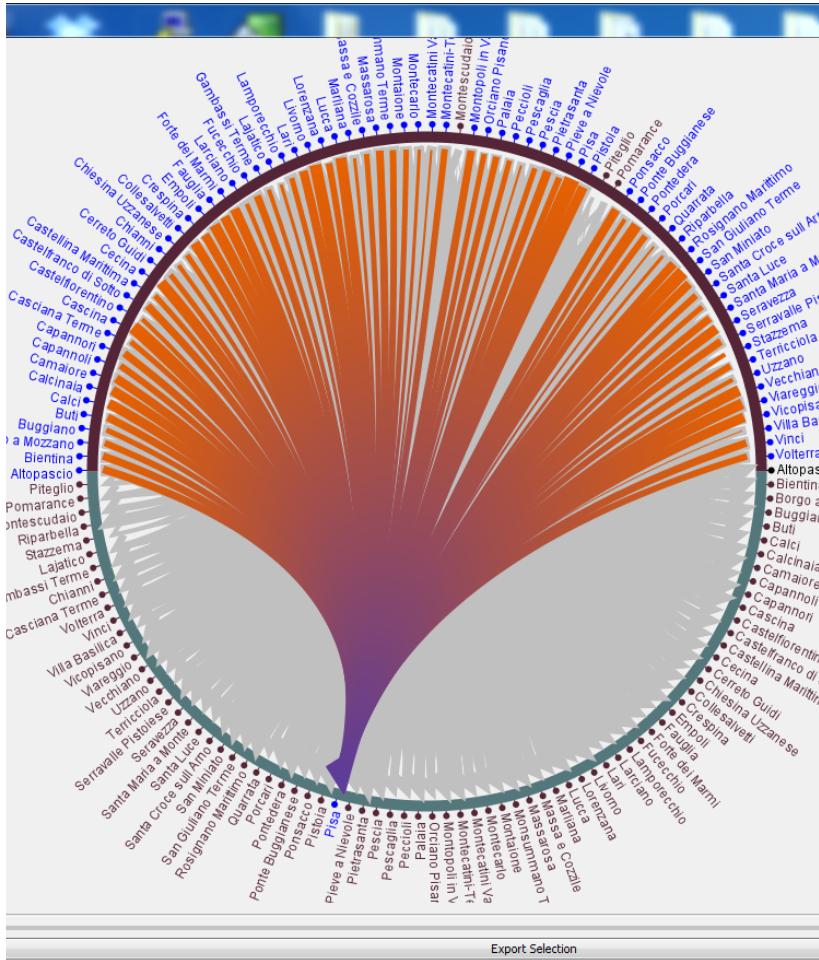# Navigate the cube at higher spatial level (provinces): visual browser



Select origins and destination from the doughnut. The map is linked with the selection. Flow weights are represented by line width
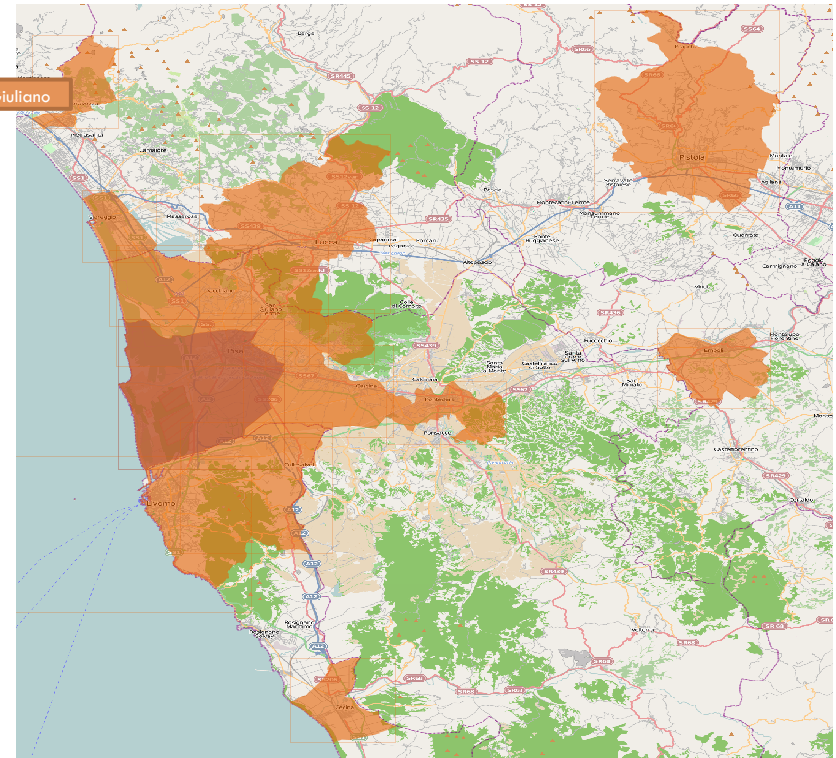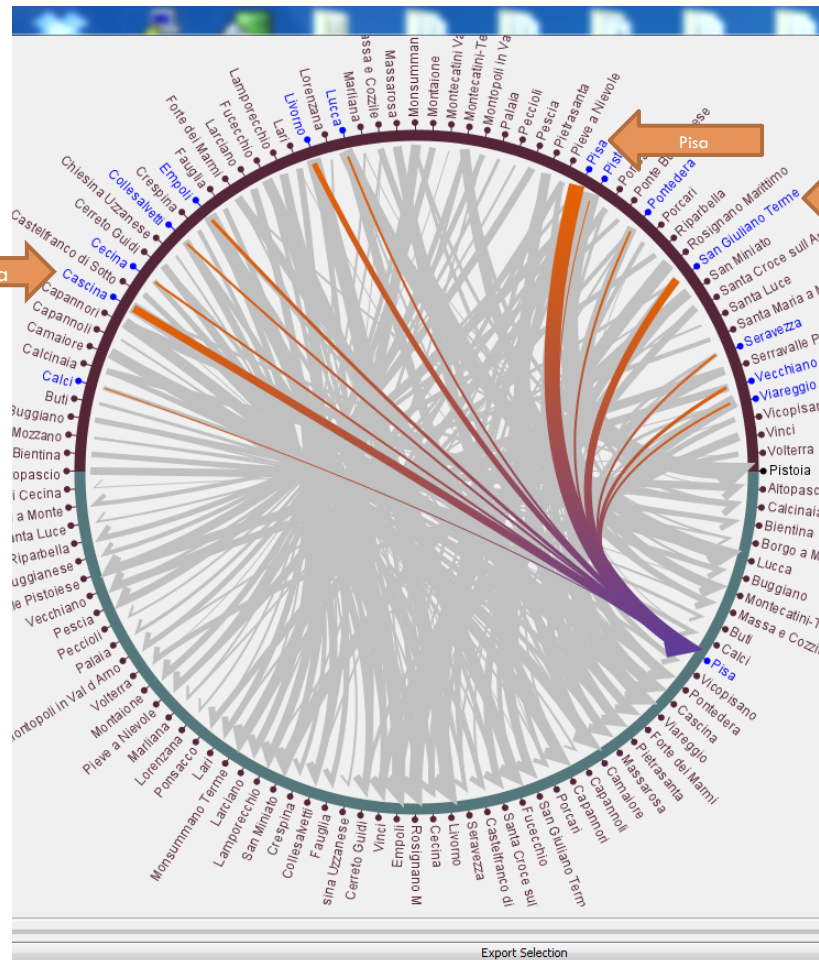
# Drill down: from province to single cities

| ime to | Cell to | | | Cell from | | Measures | |
|---|---|---|---|---|---|---|---|
| All) | Regione | Provincia | Comune | Regione | Provincia | ▼ Numero Veicoli | ◆ Pe |
| ll Time To | Toscana | ‒Pisa | | Toscana | ◆Pisa | 462.583 | 86,53 |
| | | | | | ◆Firenze | 27.742 | 13,26 |
| | | | | | ◆Livorno | 20.429 | 10,05 |
| | | | | | ◆Lucca | 17.681 | 04,07 |
| | | | | | ◆Pistoia | 5.727 | 01,30 |
| | | Pisa | ◆Pisa | Toscana | ◆Pisa | 131.430 | 28,41 |
| | | | | | ◆Livorno | 8.066 | 39,48 |
| | | | | | ◆Lucca | 7.053 | 39,89 |
| | | | | | ◆Firenze | 2.383 | 08,59 |
| | | | | | ◆Pistoia | 1.574 | 27,48 |
| | | | ◆Cascina | Toscana | ◆Pisa | 58.146 | 12,57 |
| | | | | | ◆Livorno | 1.777 | 08,70 |
| | | | | | ◆Lucca | 1.168 | 06,61 |
| | | | | | ◆Firenze | 795 | 02,87 |
| | | | | | ◆Pistoia | 305 | 05,33 |
| | | | ◆San Miniato | Toscana | ◆Pisa | 30.924 | 06,69 |
| | | | | | ◆Firenze | 12.018 | 43,32 |
| | | | | | ◆Livorno | 459 | 02,25 |
| | | | | | ◆Pistoia | 388 | 06,77 |
| | | | | | ◆Lucca | 283 | 01,60 |
| | | | ◆Pontedera | Toscana | ◆Pisa | 37.186 | 08,04 |
| | | | | | ◆Firenze | 2.402 | 08,66 |
| | | | | | ◆Livorno | 1.180 | 05,78 |
| | | | | | ◆Lucca | 611 | 03,46 |
| | | | | | ◆Pistoia | 244 | 04,26 |
| | | | ◆San Giuliano Terme | Toscana | ◆Pisa | 30.331 | 06,56 |
| | | | | | ◆Lucca | 1.983 | 11,22 |
| | | | | | ◆Livorno | 468 | 02,29 |
| | | | | | ◆Pistoia | 345 | 06,02 |
| | | | | | ◆Firenze | 124 | 00,45 |
| | | | ◆Calcinaia | Toscana | ◆Pisa | 18.425 | 03,98 |
| | | | | | ◆Livorno | 359 | 01,76 |
| | | | | | ◆Lucca | 331 | 01,87 |
| | | | | | ◆Firenze | 278 | 01,00 |
| | | | | | ◆Pistoia | 194 | 03,39 |
| | | | ◆Santa Croce sull Arno | Toscana | ◆Pisa | 14.561 | 03,15 |
| | | | | | ◆Firenze | 3.893 | 14,03 |
| | | | | | ◆Lucca | 347 | 01,96 |
| | | | | | ◆Pistoia | 290 | 05,06 |
| | | | | | ◆Livorno | 133 | 00,65 |
| | | | ◆Vecchiano | Toscana | ◆Pisa | 12.931 | 02,80 |
| | | | | | ◆Lucca | 2.700 | 15,27 |
| | | | | | ◆Pistoia | 1.032 | 18,02 |
| | | | | | ◆Livorno | 388 | 01,90 |
| | | | | | ◆Firenze | 79 | 00,28 |

□ Explode the destination by specific cities

- ■ Easy to identify the cities with the higher incomin traffic
- ■ For each city it is possible to identify the source of traffic

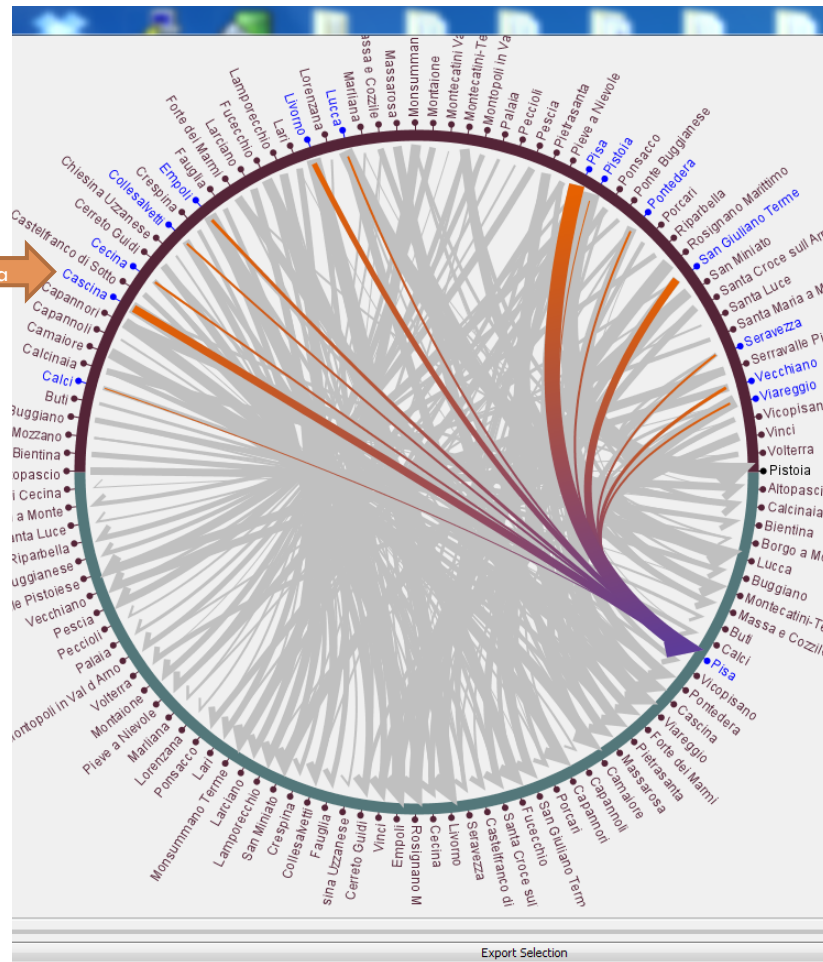# Drill down: from cities to cities

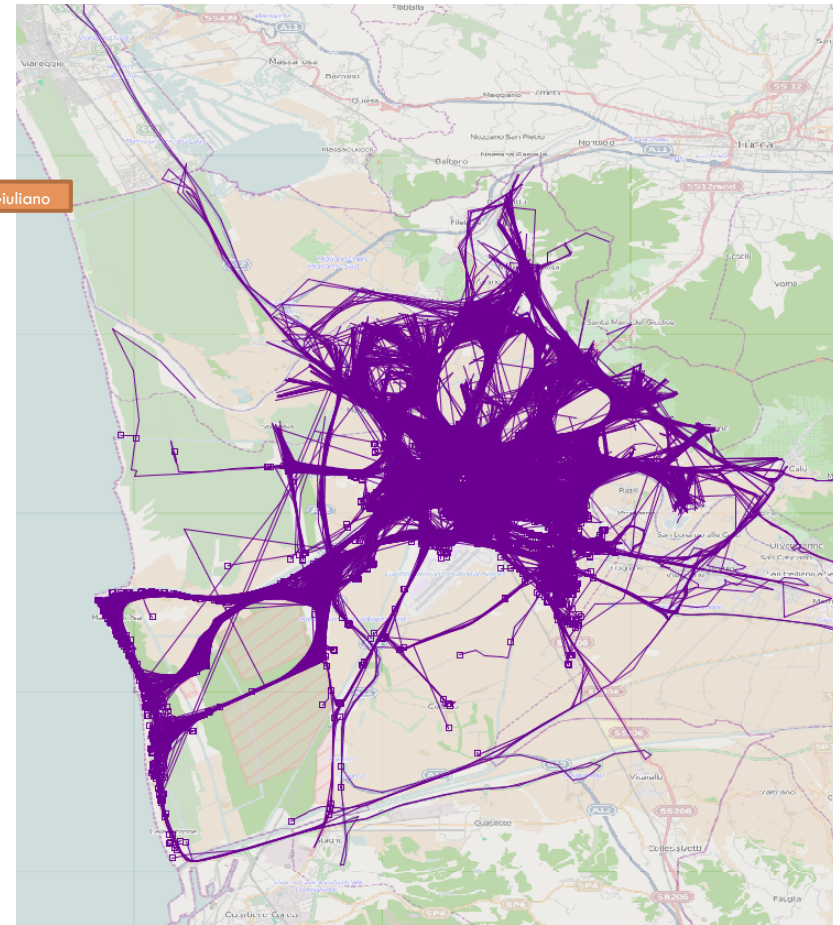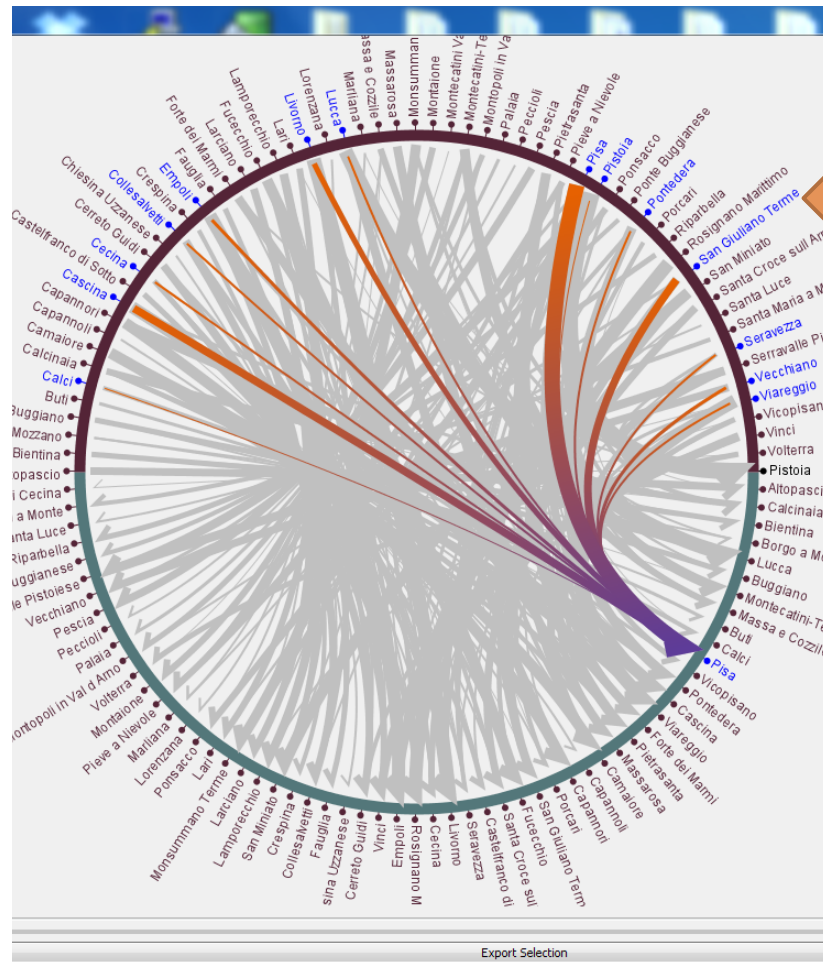# Drill down: from cities to cities (filtered)



Restrict visualization to flows above a given threshold.
Select specific flows: from Cascina, San Giuliano, and Pisa
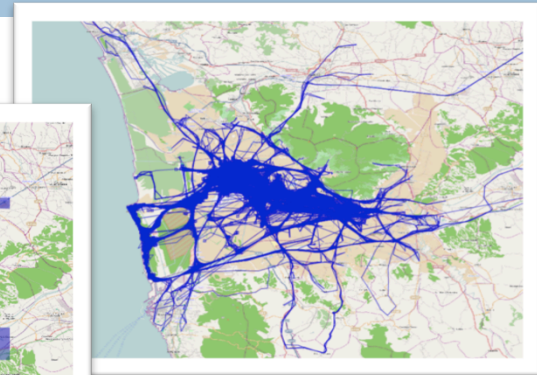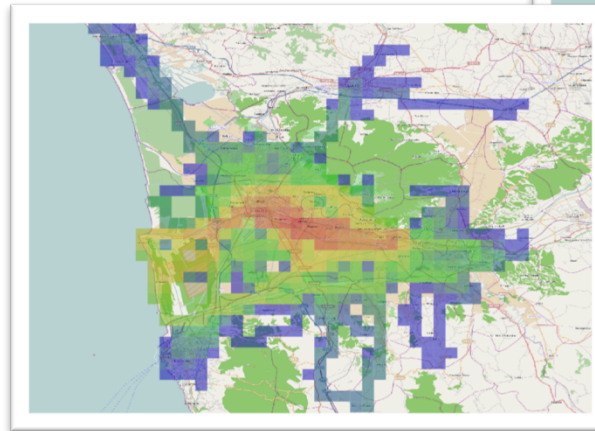
# Specific flow: from Cascina to Pisa

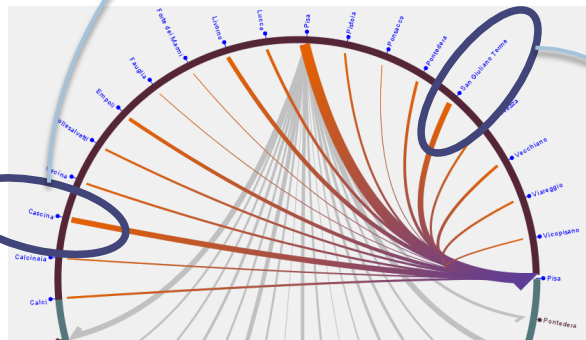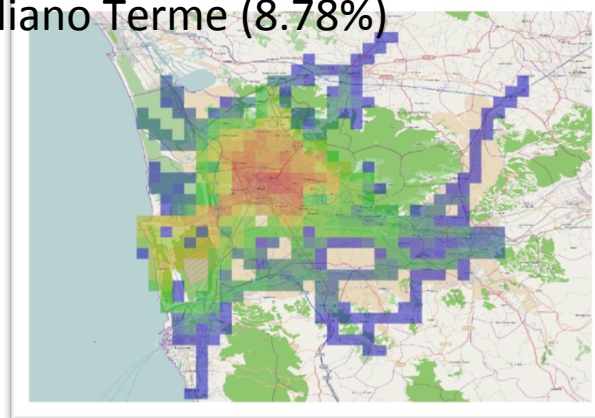# Specific flow: from San Giuliano to Pisa

# Exploring entailed data



Cascina (9.36%)
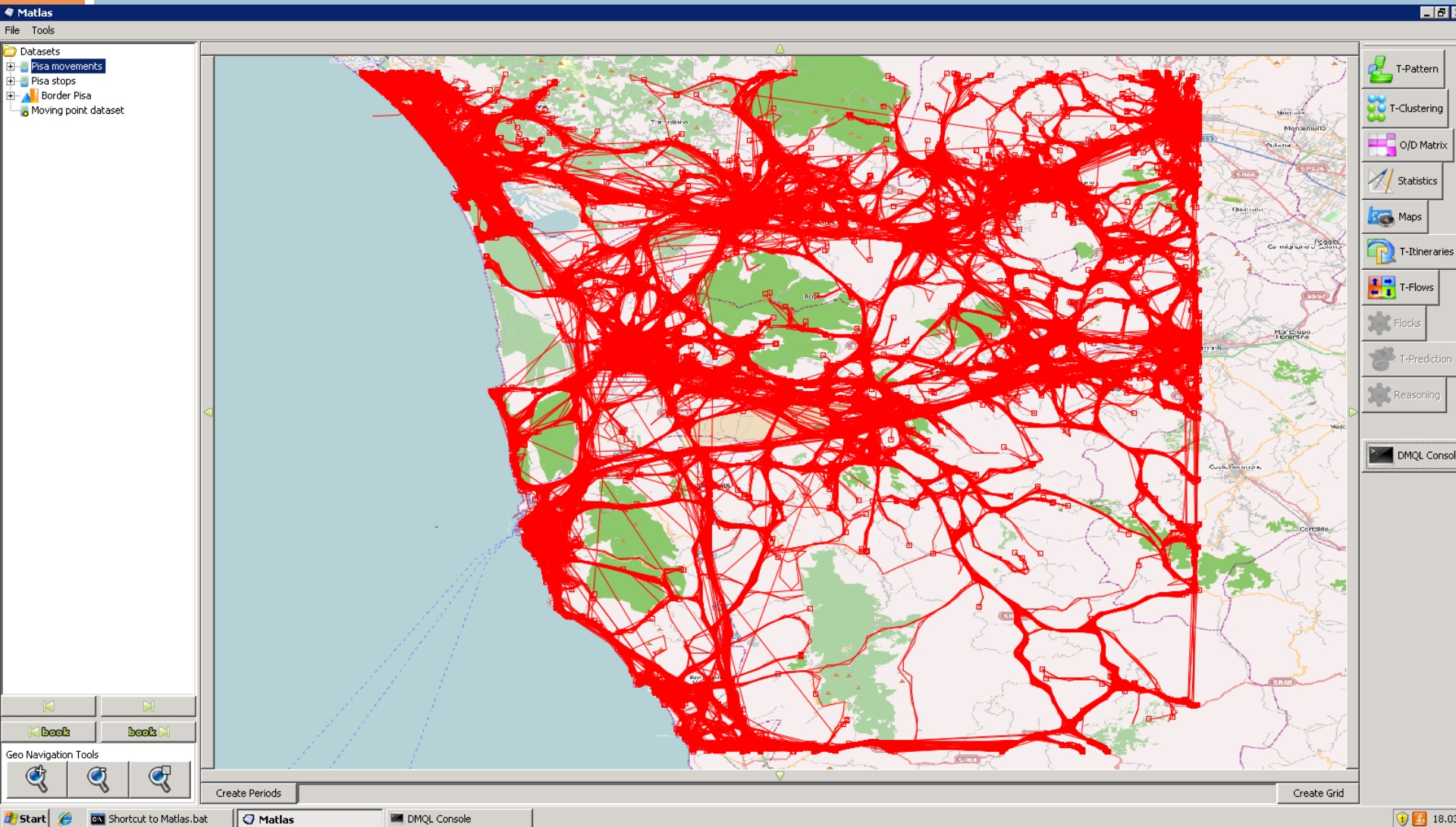
San Giuliano Terme (8.78%)

| Cell To | Cell From | Measures | |
|---|---|---|---|
| | | Num vehicles | % |
| +Pisa | −Pisa | 89.730 | 84,24% |
| | +Pisa | 63.331 | 70,58% |
| | +Cascina | 8.402 | 09,36% |
| | +San Giuliano Terme | 7.877 | 08,78% |
| | +Vecchiano | 1.869 | 02,08% |
| | +Pontedera | 1.408 | 01,57% |
| | +Calci | 1.220 | 01,36% |

Discovering **access patterns** to Pisa with GPS track data

# **Access patterns** using T-clustering

# Access patterns using T-clustering

# **Access patterns** using T-clustering



*Marina di Pisa/Tirrenia*

*Lucca*

*A12 Sud*

*Cascina*

# Characterizing the **access patterns**: origin & time



1,50 %

*A12 Sud*

Origin distribution



**Distribuzione Origini**

- Pisa
- Marina/Tirrenia
- A12 (Nord)
- FiPiLi (Empoli)
- A12 (Sud)
- Lucca
- A11 (Pistoia)
- Collesalvetti
- Ponsacco
- SS12 (Nord Lucca)
- Montecatini
- Torre del Lago
- Calci
- Asciano
- Altre origini
  Rumore

2,90 %

*Marina di Pisa/Tirrenia*

# Persistency of **access patterns**

# Studying the attractiveness/efficiency of a service with GPS tracks
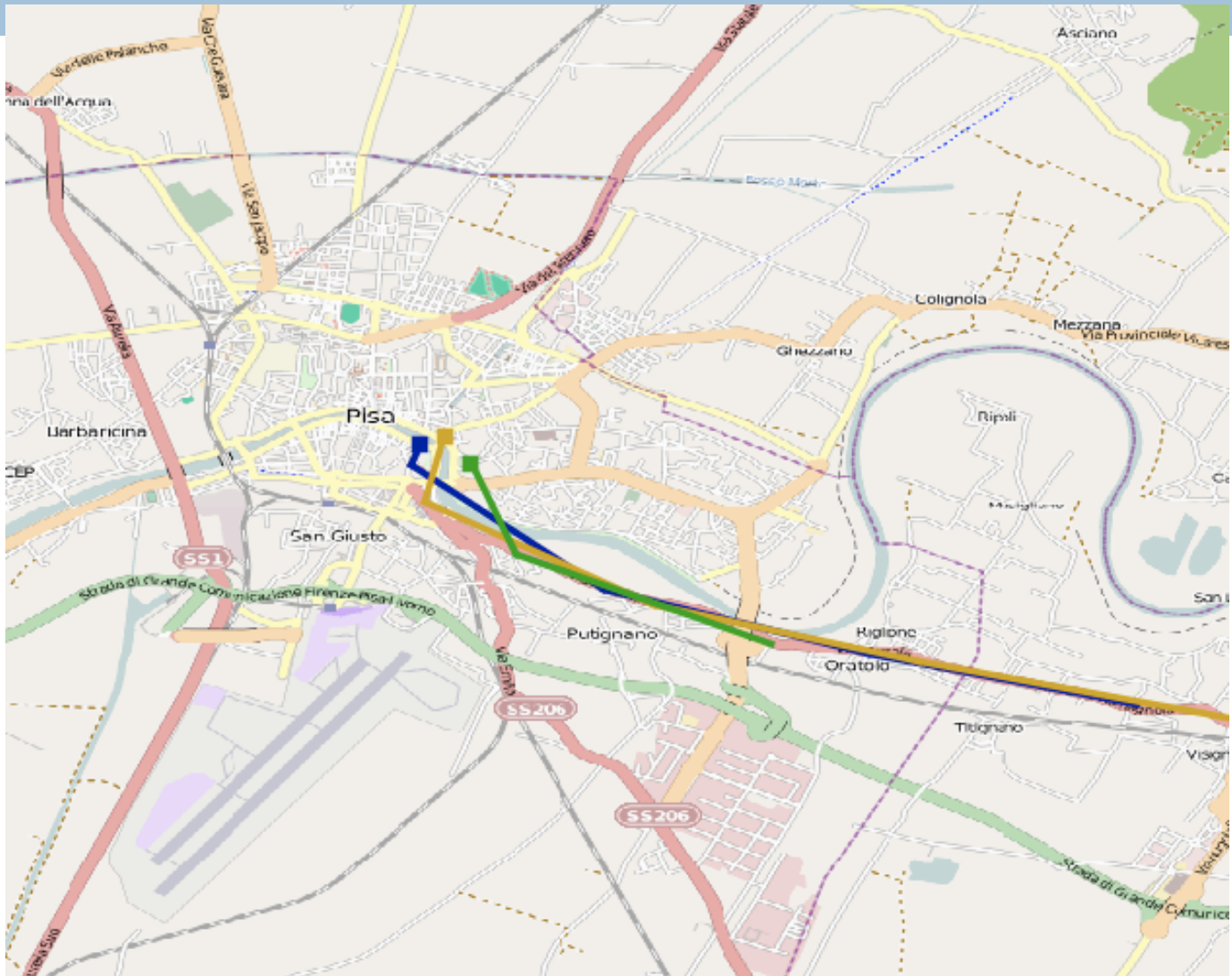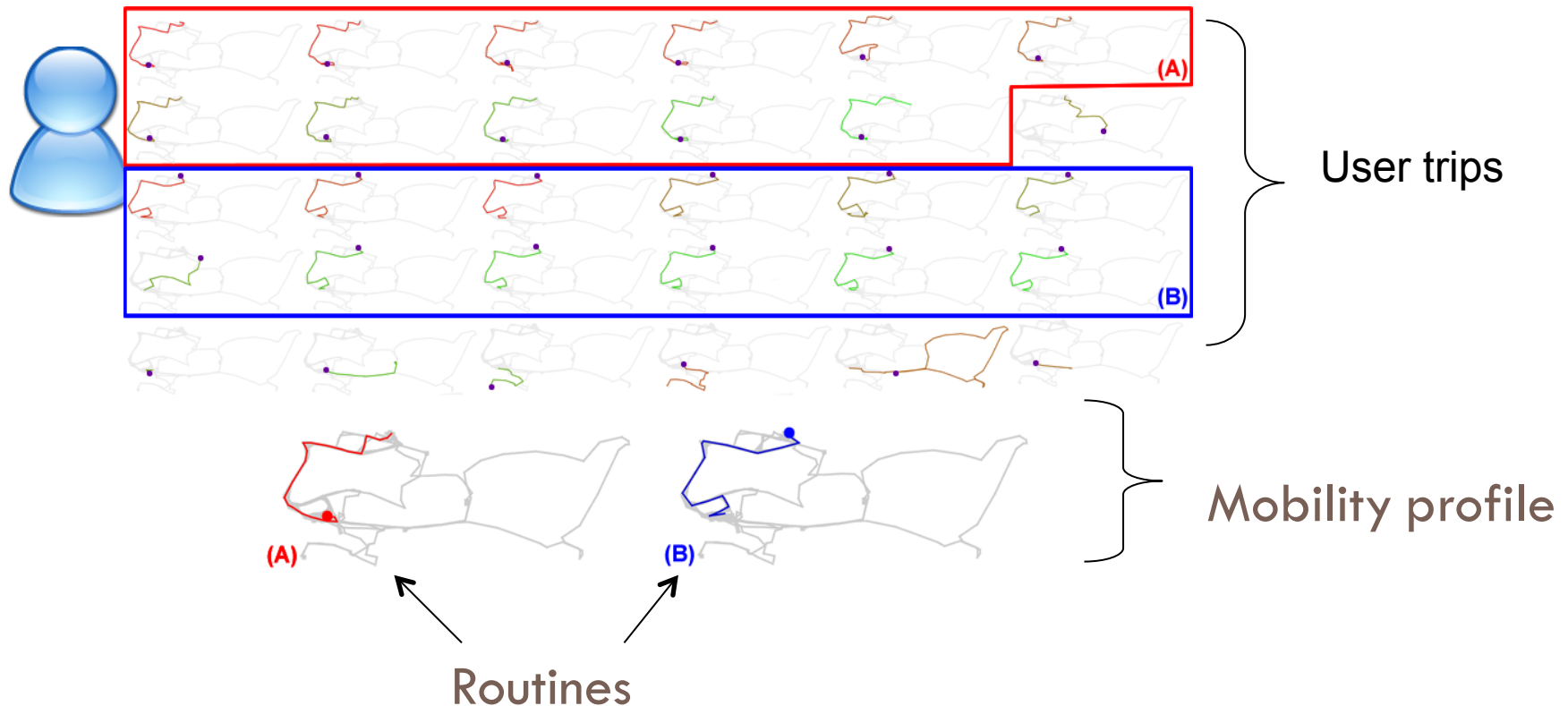
# Discovering mobility profiles with GPS tracks data

# Extract travellers profiles

# Extracting travellers profiles

- Analysis focused on the single individual

- Find his/her systematic mobility



User trips

Mobility profile

Routines

# Application: Car pooling

Pro-active suggestions of sharing rides opportunities without the need for the user to explicitly specify the trips of interest.
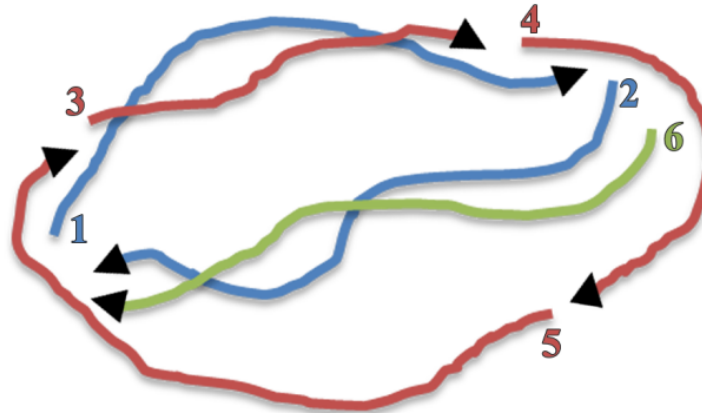
Matching two routines:

$$contained(T_1, T_2, th_{distance}^{walking}, th_{time}^{wasting}) \equiv \exists i, j \in \mathcal{N} \mid$$
$$0 < i \leq j \leq m \wedge$$
$$Dist(p_1^1, p_i^2) + Dist(p_n^1, p_j^2) \leq th_{distance}^{walking} \wedge$$
$$Dur(p_1^1, p_i^2) + Dur(p_n^1, p_j^2) \leq th_{time}^{wasting}$$

Mobility profile share-ability:

$$mobility\ profiles\ \tilde{T}_1\ and\ \tilde{T}_2$$

$$profileShare(\tilde{T}_1, \tilde{T}_2, th_{distance}^{walking}, th_{time}^{wasting}) =$$

$$\frac{\left| \left\{ p \in \tilde{T}_1 \mid \exists q \in \tilde{T}_2 . Share(p, q, th_{distance}^{walking}, th_{time}^{wasting}) \right\} \right|}{\mid \tilde{T}_1 \mid}$$



|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | - | - | F | F | F | F |
| 2 | - | - | F | F | F | T |
| 3 | T | F | - | - | - | F |
| 4 | F | F | - | - | - | F |
| 5 | F | F | - | - | - | F |
| 6 | F | T | F | F | F | - |

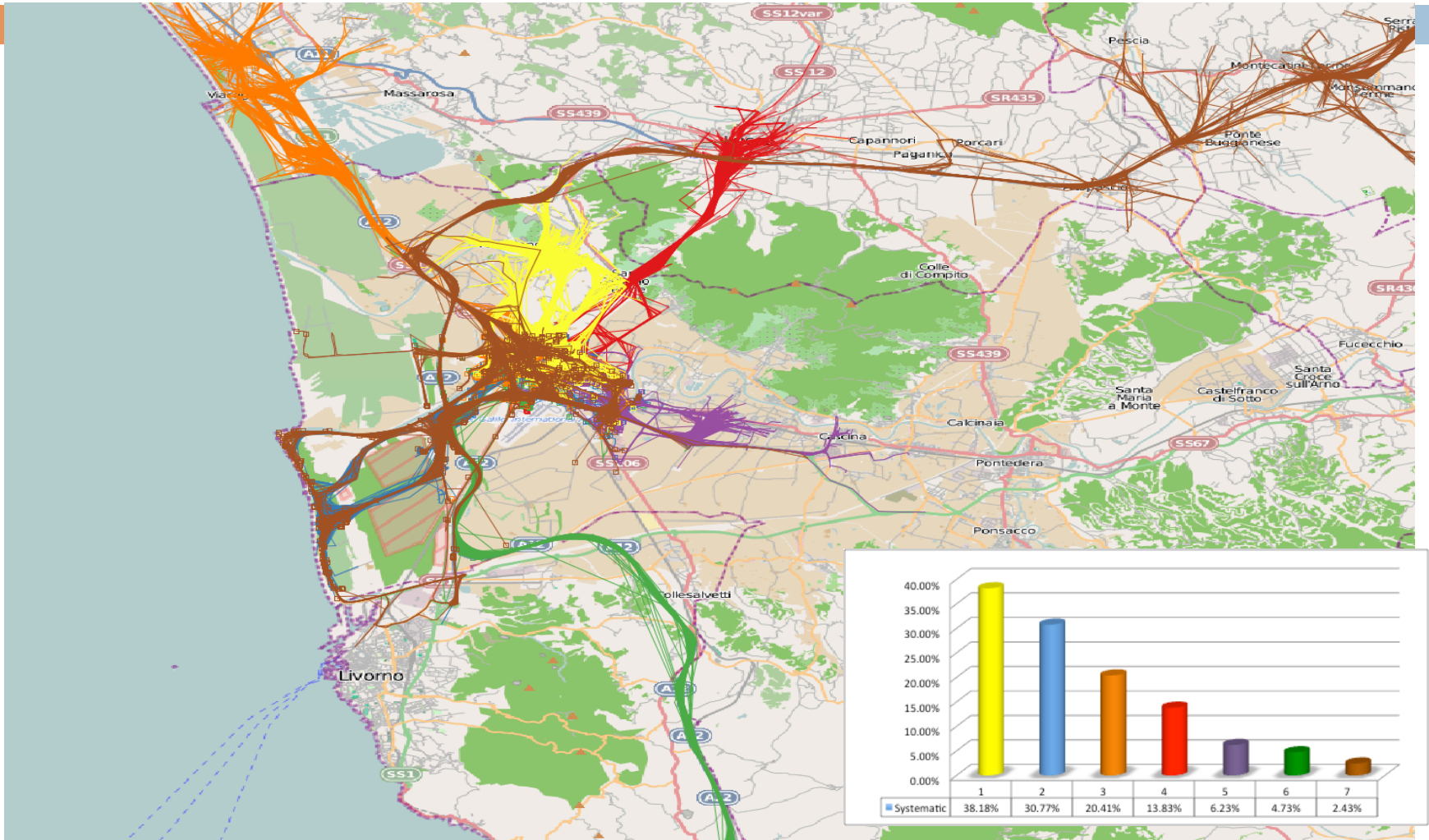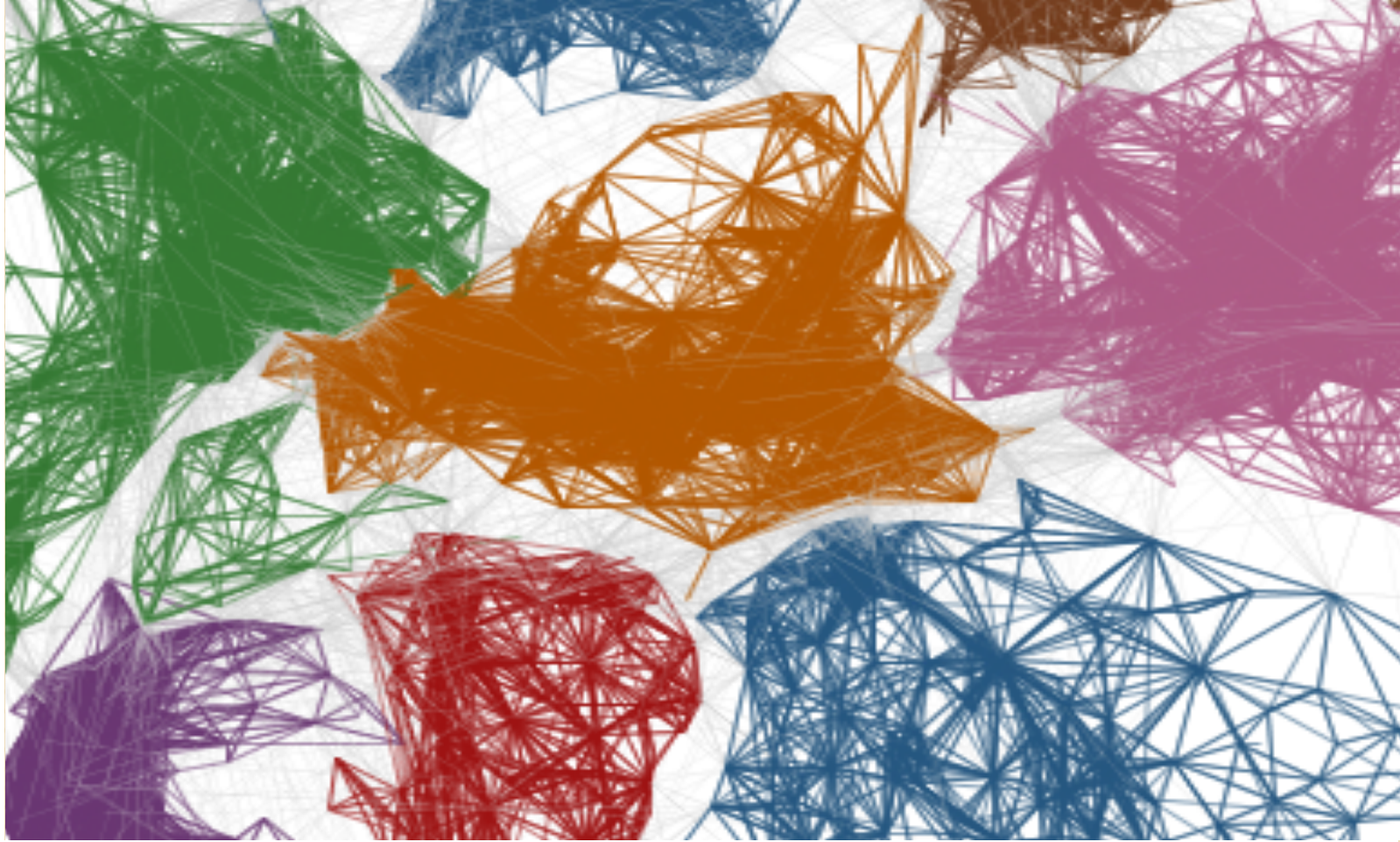|   |   |   |   |
|---|---|---|---|
|   | - | 0 | 1/2 |
|   | 1/3 | - | 0 |
|   | 1 | 0 | - |

# Car pooling potential

67.2% routines match with a routine of other users

32.5% users share one or more routines with other users

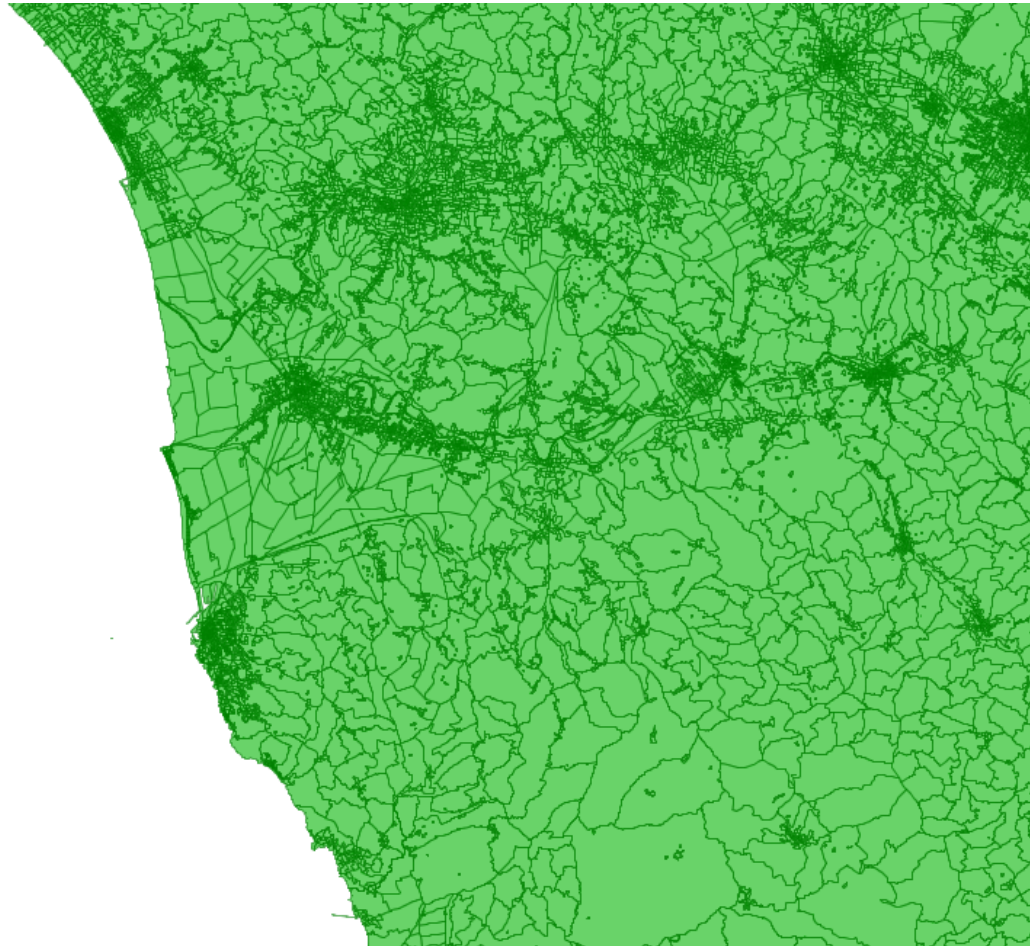# Impact of systematic mobility on access patterns
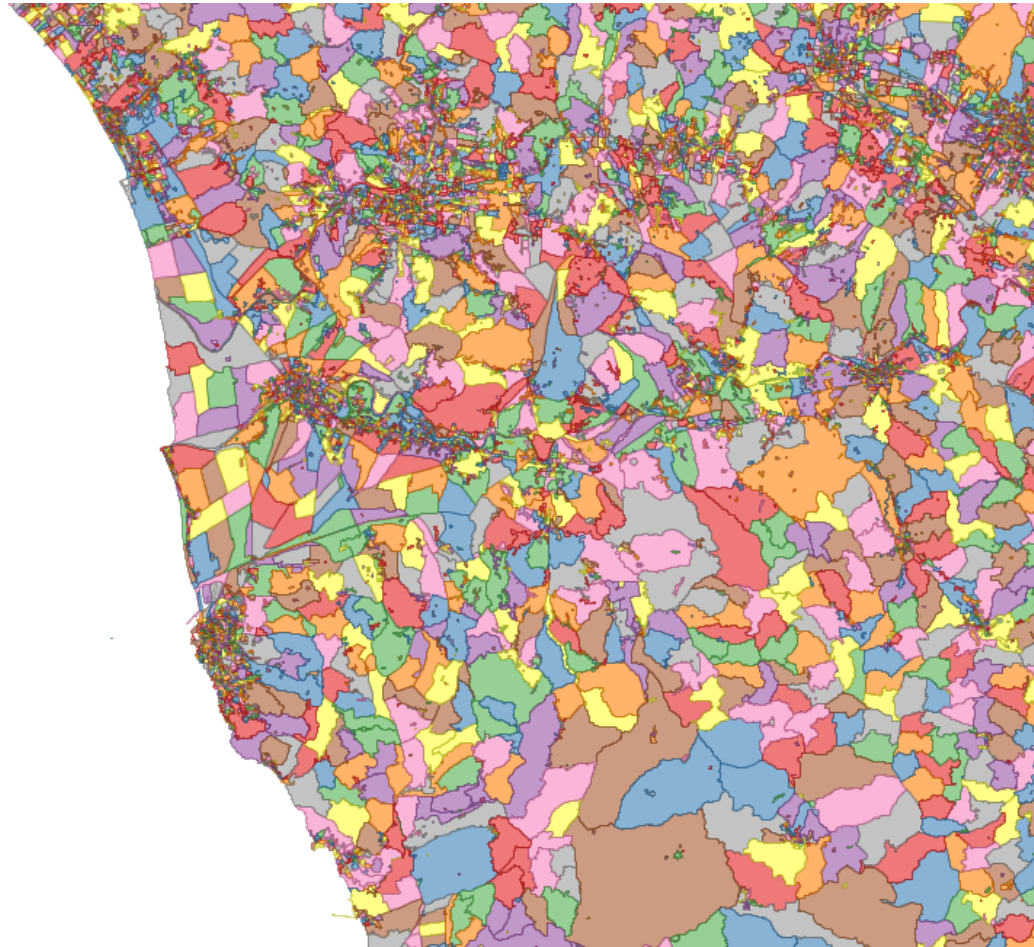
Find border of human mobility

# Motivations

- Mobility management offices needs accurate information to handle mobility issues
  - Monitoring: how to predict/manage emergences of special events?
  - Planning: public transportation desisgn, incentives for multi-modal movement, etc.
- Planning involse several entities
  - The city level is not sufficient: the neighbor cities are necessarily influenced
  - The regional level is too general: lost focus for specific/local requirements
  - Does provinces provide the necessary level of details?

# Step 1: spatial regions
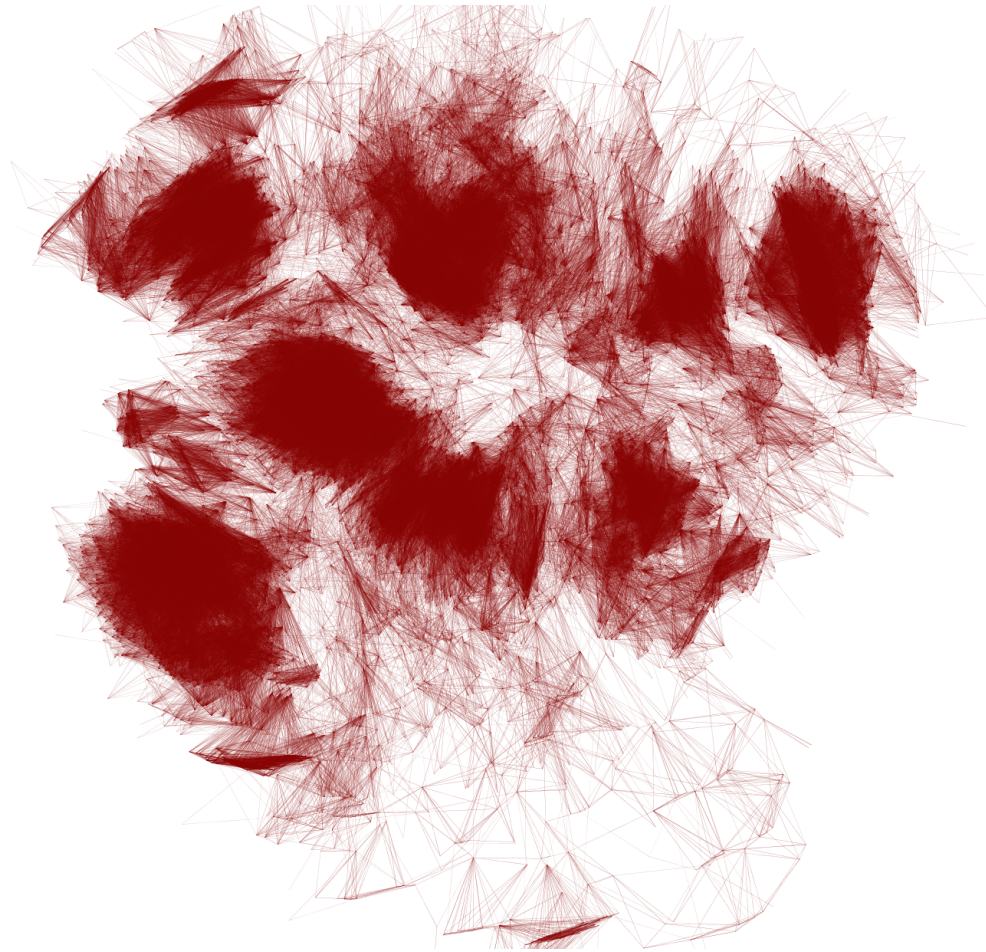
# Start from random labeling for region

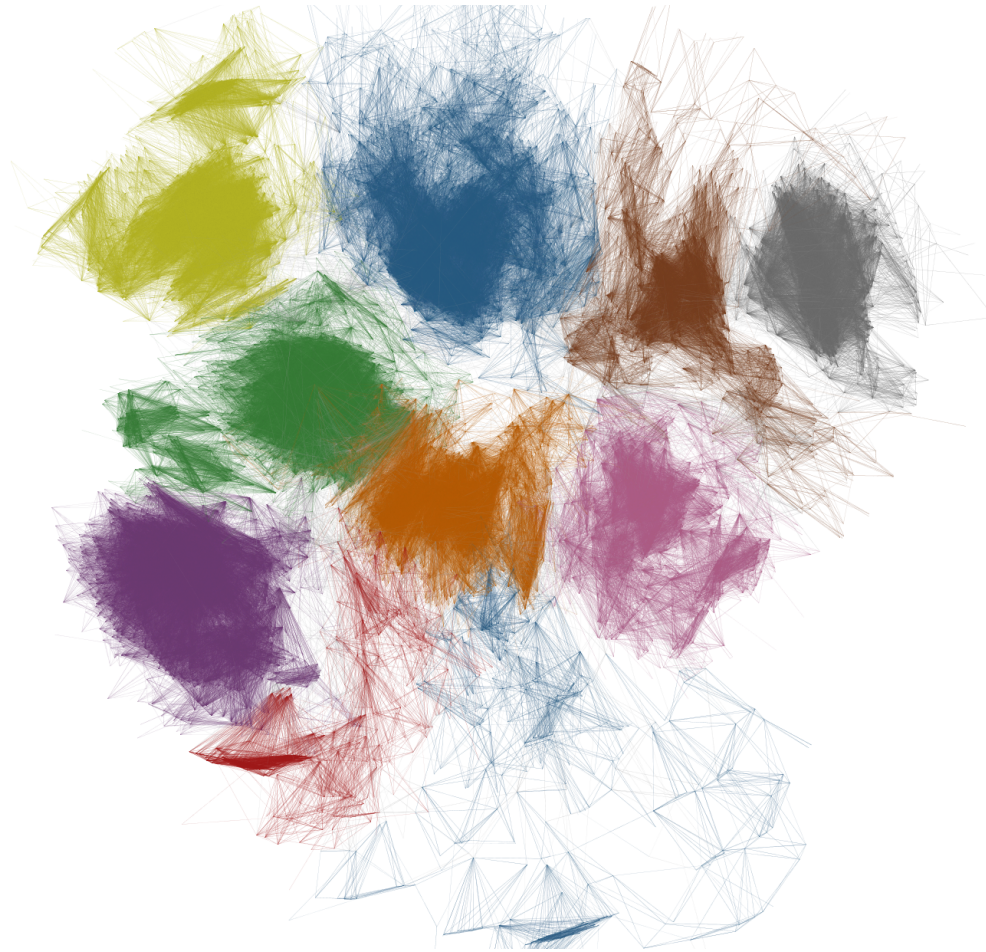# Step 2: evaluate flows among regions
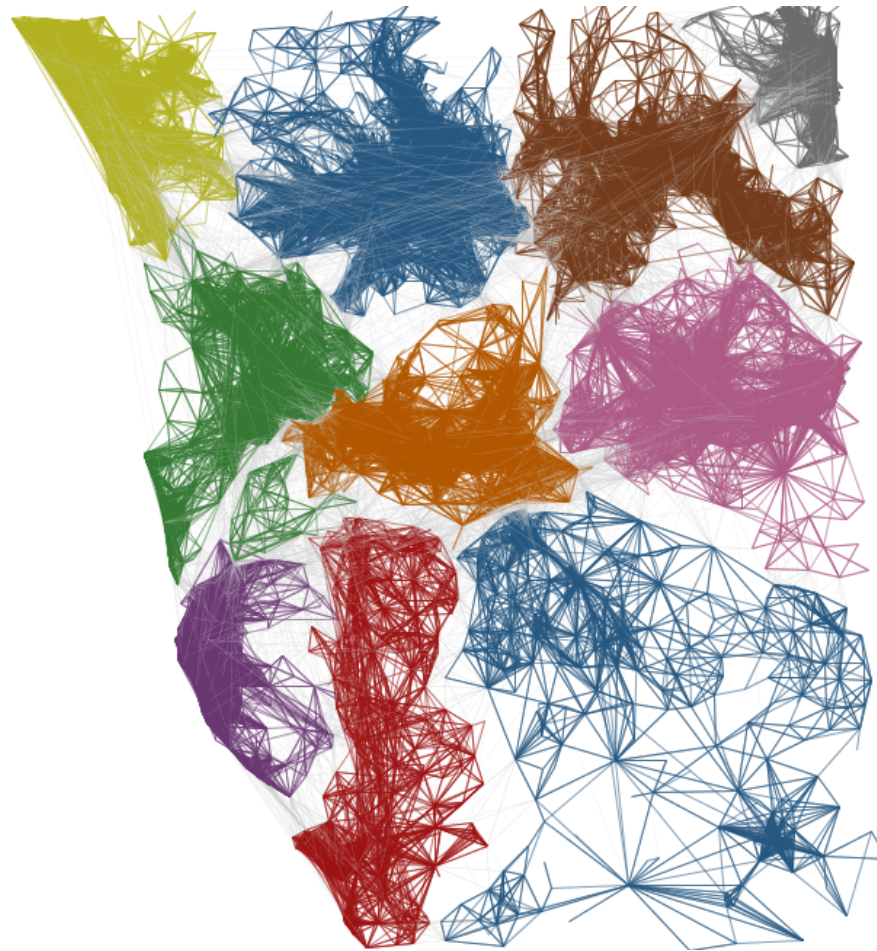
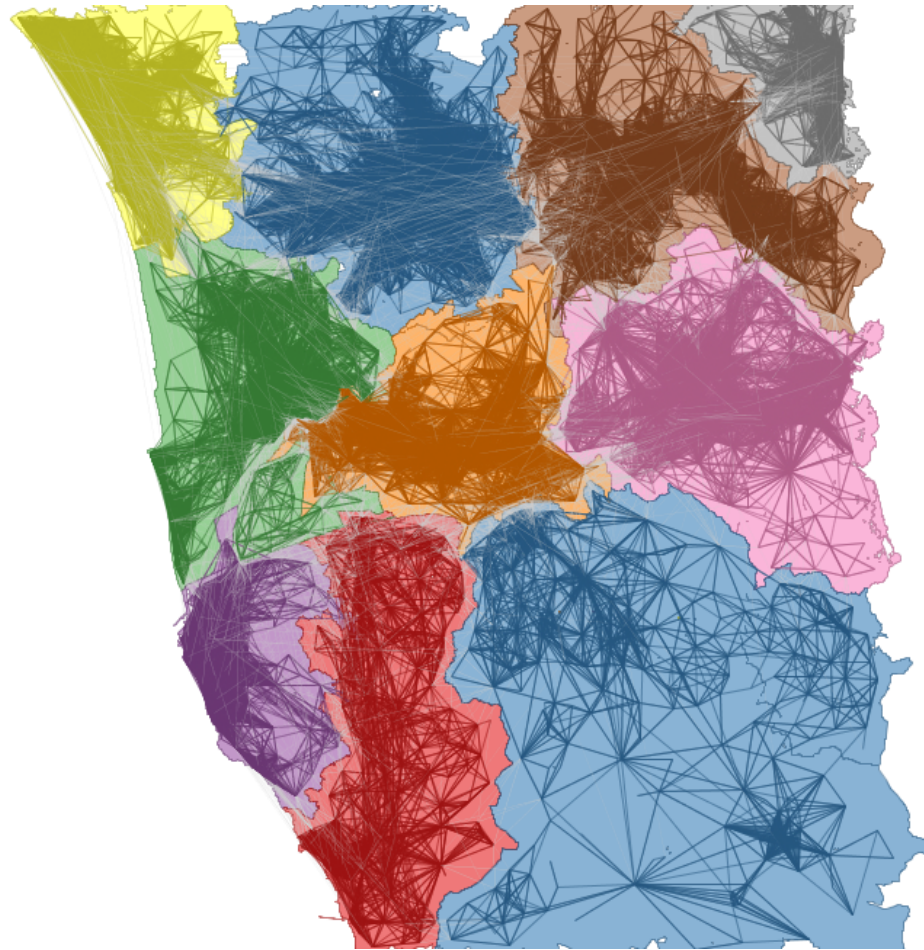# Step 3: consider only the network

# Step 4: perform clustering
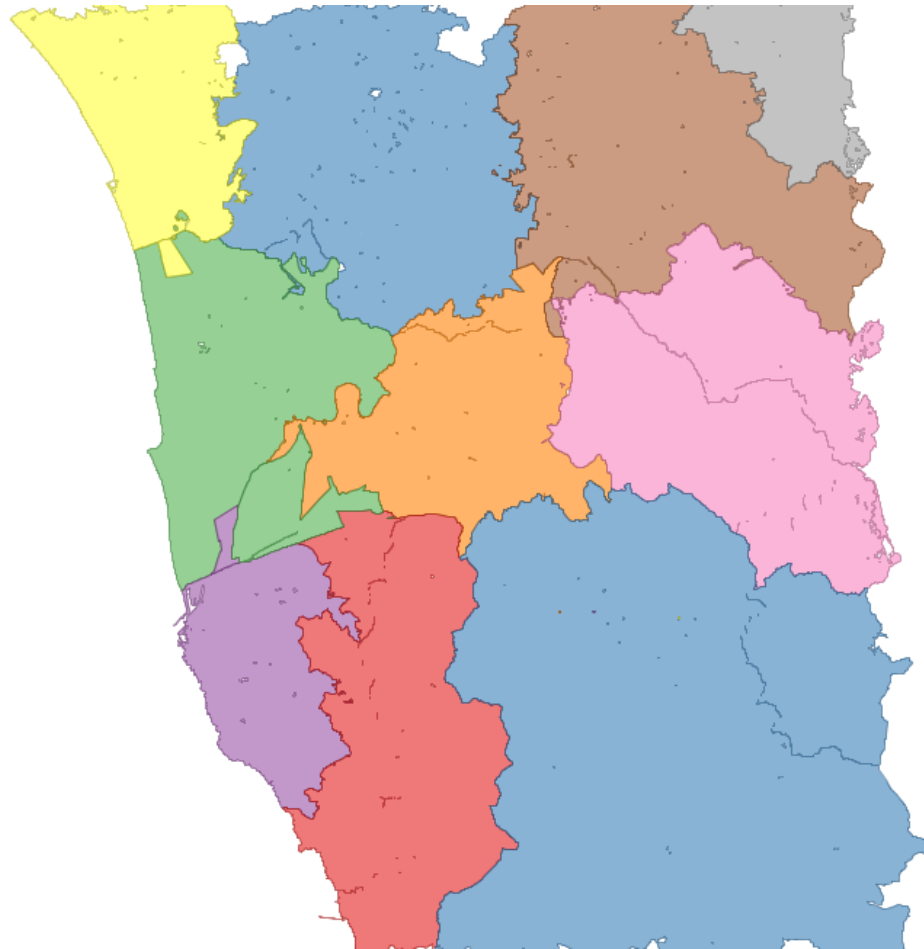
# Step 4: perform clustering
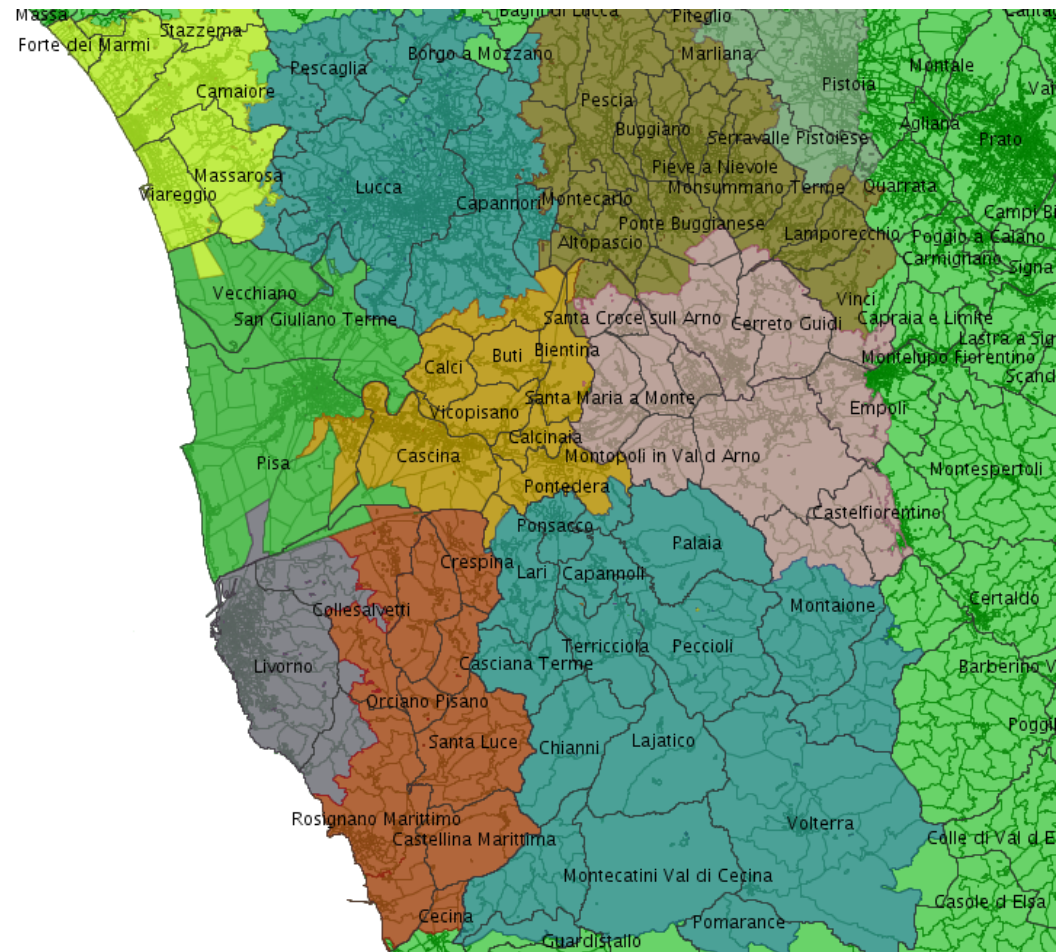
# Step 5: map nodes back to geography

# Step 5: map nodes back to geography

# Final result

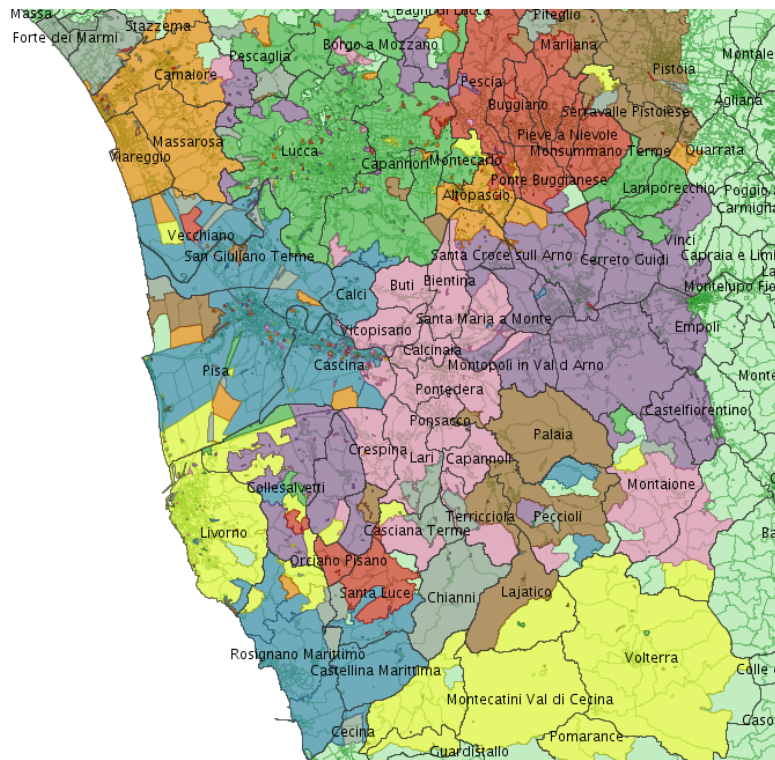# Final result: comparison
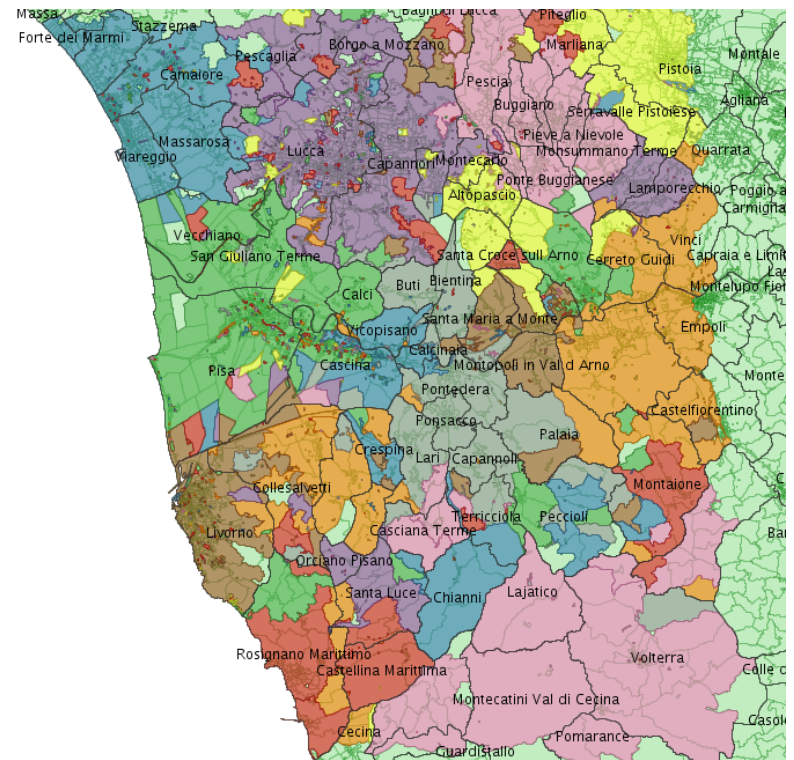
# Borders using only OD flows

# Borders in different time periods



**Only weekdays movements**



Similar to global clustering: strong influence of systematic movements
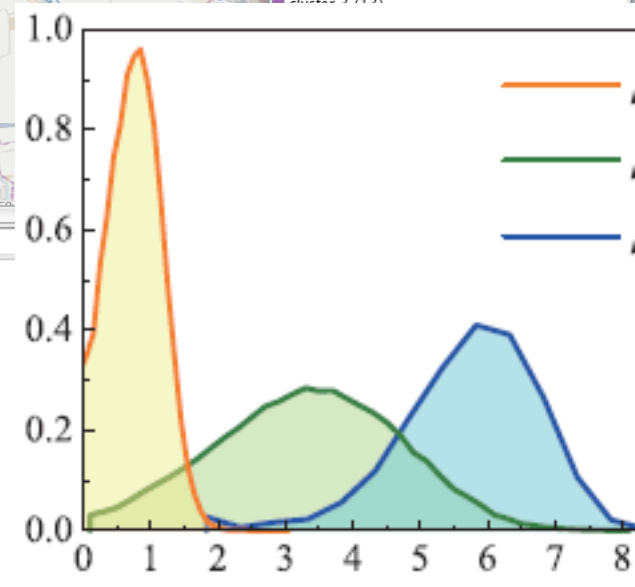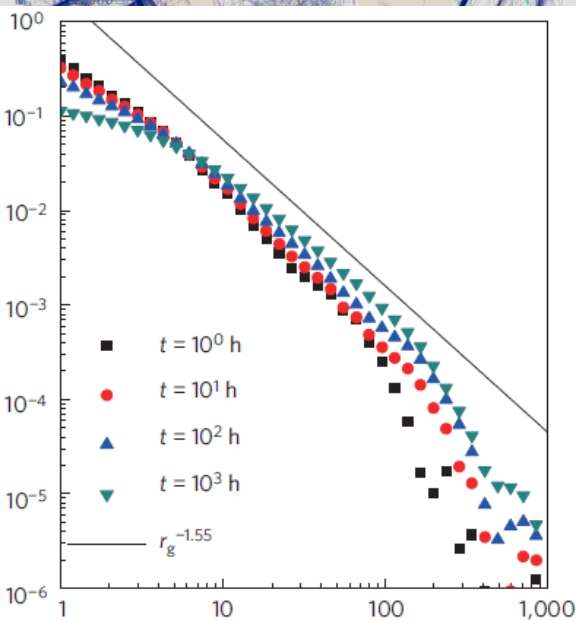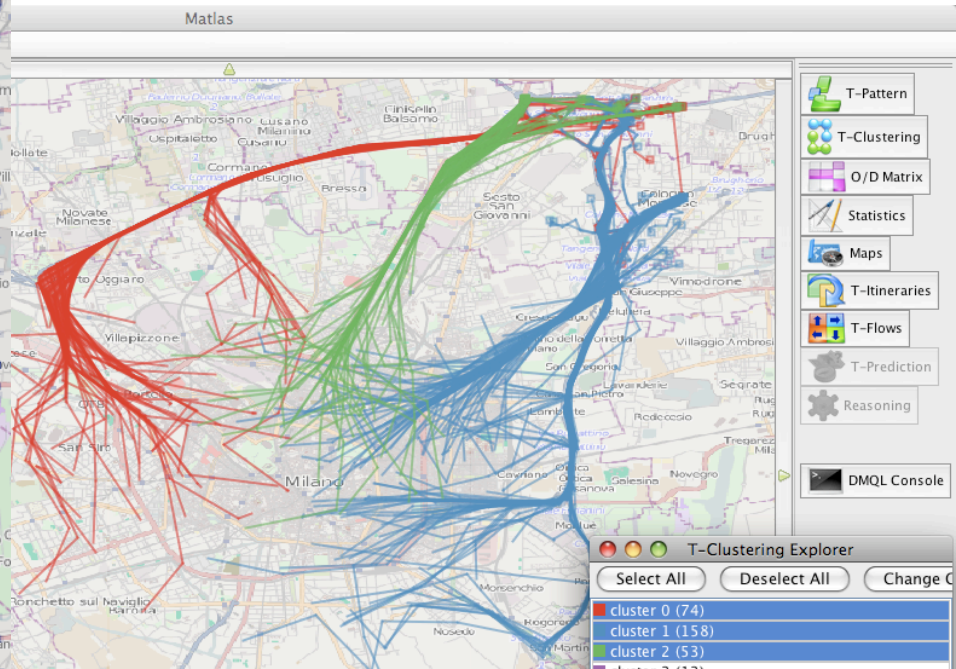
**Only weekend movements**

Strong fragmentation: the influence of systematic movements (home-work) is missing

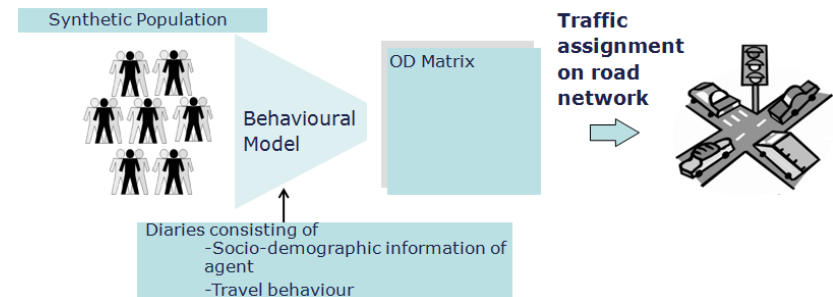# Summarizing: big data push towards converging sciences

# Big data push towards convergence

- Network science
  - **Global models** of complex social phenomena
  - Behavioral **diversity** in society at large
- Data mining
  - **Local patterns** of complex social phenomena
  - Behavioral **similarity** in sub-populations
- Both visions needed to achieve realistic and accurate models for prediction and **simulation**
  - Computational sociology (Lazer et al., Science 2009)
- Both data-driven, each leverage on the other

# DATA-SIM – Data science for simulating the era of electric vehicles

- What's the impact on mobility and energy distribution in the case of a massive switch to electric cars?

- Data mining + network science + agent-based simulation

- FET project started October 2011 www.datasimfet.eu

- KDD LAB Pisa + I-MOB Hasselt + Barabasi Lab Budapest+OCTO



Synthetic Population

Behavioural Model

OD Matrix

Traffic assignment on road network

Diaries consisting of
-Socio-demographic information of agent
-Travel behaviour

# Knowledge Discovery and Data Mining Laboratory

**Web Site: http://kdd.isti.cnr.it**

*Personnel*

**Lab Head**



Cappelli Amedeo

Giannotti Fosca

Pedreschi Dino

Turini Franco

**Post Doc**

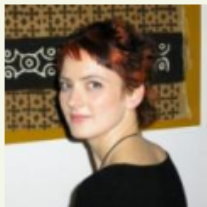Berlingerio Michele

Pinelli Fabio
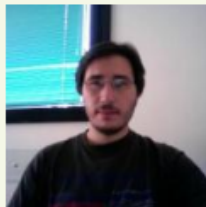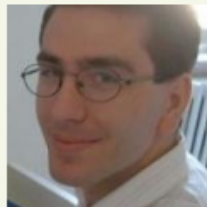
Trasarti Roberto

**Research Staff**

Nanni Mirco

Renso Chiara

Rinzivillo Salvatore

Ruggieri Salvatore

**PhD Student**

Coscia Michele

Monreale Anna

Ong Rebecca

Pennacchioli Diego

Caterina D'angelo

Claudio Schifani

Chiara Falchi

Zehui Qu

Barbara Furletti, Andrea Romei, Sergio Barsocchi